# 6.5 Data Compression

‣ basics
‣ run-length encoding
‣ Huffman compression
‣ LZW compression

*References: Algorithms 4th edition, Section 6.5*

---

## 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.

---

## 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.

---

## Lossless compression and expansion

**Message.** Binary data $B$ we want to compress.   ← uses fewer bits (you hope)
**Compress.** Generates a "compressed" representation $C(B)$.
**Expand.** Reconstructs original bitstream $B$.

| bitstream $B$ | **Compress** | compressed version $C(B)$ | **Expand** | original bitstream $B$ |
|---|---|---|---|---|
| 0110110101... | | 1101011111... | | 0110110101... |

**Basic model for data compression**

**Compression ratio.** Bits in $C(B)$ / bits in $B$.

**Ex.** 50-75% or better compression ratio for natural language.

## 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?

---

‣ **binary I/O**
‣ genomic encoding
‣ run-length encoding
‣ Huffman compression
‣ LZW compression

---

## Reading and writing binary data

Binary standard input and standard output.   Libraries to read and write bits from standard input and to standard output.

```
public class BinaryStdIn
    boolean readBoolean()        read 1 bit of data and return as a boolean value
       char readChar()           read 8 bits of data and return as a char value
       char readChar(int r)      read r bits of data and return as a char value
    [similar methods for byte (8 bits); short (16 bits); int (32 bits); long and double (64 bits)]
    boolean isEmpty()            is the bitstream empty?
       void close()              close the bitstream
```
**API for static methods that read from a bitstream on standard input**

```
public class BinaryStdOut
       void write(boolean b)     write the specified bit
       void write(char c)        write the specified 8-bit char
       void write(char c, int r) write the r least significant bits of the specified char
    [similar methods for byte (8 bits); short (16 bits); int (32 bits); long and double (64 bits)]
       void close()              close the bitstream
```
**API for static methods that write to a bitstream on standard output**

---

## Reading and writing binary data

Date representation.  Consider different ways to represent 12/31/1999.

**A character stream (StdOut)**
```
StdOut.print(month + "/" + day + "/" + year);
```
00110001 00110010 00101111 00110011 00110001 00101111 00110001 00111001 00111001 00111001
1        2        /        3        1        /        1        9        9        9      80 bits

**Three ints (BinaryStdOut)**
```
BinaryStdOut.write(day);
BinaryStdOut.write(month);
BinaryStdOut.write(year);
```
00000000000000000000000000001100 00000000000000000000000000011111 00000000000000000000011111001111
12                               31                               1999                            96 bits

**Two chars and a short (BinaryStdOut)**
```
BinaryStdOut.write((char) day);
BinaryStdOut.write((char) month);
BinaryStdOut.write((short) year);
```
00001100 00011111 0000011111001111
12       31       1999            32 bits

**A 4-bit field, a 5-bit field, and a 12-bit field (BinaryStdOut)**
```
BinaryStdOut.write(day, 4);
BinaryStdOut.write(month, 5);
BinaryStdOut.write(year, 12);
```
1100 11111 011111001111 000
12   31    1999          21 bits ( + 3 bits for byte alignment at close)

## Binary dumps

**Q.** How to examine the contents of a bitstream?



Hexadecimal to ASCII conversion table

**Standard character stream**
```
% more abra.txt
ABRACADABRA!
```

**Bitstream represented as 0 and 1 characters**
```
% java BinaryDump 16 < abra.txt
0100000101000010
0101001001000001
0100001101000001
0100010001000001
0100001001010010
0100000100100001
96 bits
```

**Bitstream represented with hex digits**
```
% java HexDump 4 < abra.txt
41 42 52 41
43 41 44 41
42 52 41 21
96 bits
```

**Bitstream represented as pixels in a Picture**
```
% java PictureDump 16 < abra.txt
```

*16-by-6 pixel window, magnified*
```
96 bits
```

**Four ways to look at a bitstream**

---

‣ binary I/O
‣ **limitations**
‣ genomic encoding
‣ run-length encoding
‣ Huffman compression
‣ LZW compression

---

## Universal data compression

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.… "*

---

## 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

**Proposition.** No algorithm can compress every bitstream.

**Pf 2.** [by contradiction]
- Suppose you have a universal data compression algorithm U that can compress every bitstream.
- Given bitstream $B_0$, compress it to get smaller bitstream $B_1$.
- Compress $B_1$ to get a smaller bitstream $B_2$.
- Continue until reaching bitstream of size 0.
- Implication: all bitstreams can be compressed with 0 bits!

**Pf 1.** [by counting]
- Suppose you have a compression algorithm that can compress all 1,000-bit streams.
- $2^{1000}$ possible bitstreams with 1000 bits.
- Only $1 + 2 + 4 + \ldots + 2^{998} + 2^{999}$ can be encoded with $\leq 999$ bits.
- Similarly, only 1 in $2^{499}$ bitstreams can be encoded with $\leq 500$ bits!

Universal
data compression?

13

```
% java RandomBits | java PictureDump 2000 500
```

1000000 bits

**A difficult file to compress: one million (pseudo-) random bits**

```java
public class RandomBits
{
    public static void main(String[] args)
    {
        int x = 11111;
        for (int i = 0; i < 1000000; i++)
        {
            x = x * 314159 + 218281;
            BinaryStdOut.write(x > 0);
        }
        BinaryStdOut.close();
    }
}
```

14

| | | | |
|---|---|---|---|
| absolutely essential | convicted felon | inner core | protective helmet |
| ACT test | crimson red | jet plane | raise up |
| advance forward | delete out | join together | repeat again |
| advance warning | down under | KFC chicken | revert back |
| at this point in time | Each and every one of you | kitty cat | rising above |
| attach together | elevate up | last will and testament | safe sanctuary |
| awful bad | end result | LCD display | Scotch Whisky |
| bad trouble | enter into | live audience | sharp point |
| bare naked | evil villain | male son | sink down |
| basic fundamentals | falling down | marital spouse | small speck |
| begin to proceed | famous celebrities | merge together | solitary hermit |
| boiling hot | fellow colleagues | new discovery | specific example |
| bunny rabbit | first of all | newborn baby | square box |
| cease and desist | for your FYI | null and void | still remains |
| Chile pepper | foreign imports | original founder | swampy marsh |
| circulated around | free gift | over and above | temper tantrum |
| close proximity | full satisfaction | PIN number | terrible tragedy |
| close scrutiny | gather together | pair of twins | tiny bit |
| closed fist | grand total | past experience | true fact |
| collaborate together with | handwritten manuscript | past tradition | turning around |
| combined together | HIV virus | positive yes | under cover |
| complete monopoly | hopes and aspirations | postponed until later | unique individual |
| completely filled | hygienic cleaning | previous history | useless and unnecessary |
| component parts | immortalized forever | print out | wall murals |
| continuing on | individual person | proposed plan | whether or not |

http://www.corsinet.com/braincandy/twice.html

15

**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 denmtrasote. In a pubiltacion of New Scnieitst you could ramdinose all the letetrs, keipeng the first two and last two the same, and reibadailty 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 sugsegts we may have some pofrweul palrlael 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 meniang. "* — *Graham Rawlinson*

**A.** Quite a bit.

16

## Slide 17

‣ **genomic encoding**
‣ run-length encoding
‣ Huffman compression
‣ LZW compression

---

## Slide 18

### Genomic code

Genome. String over the alphabet { A, C, T, G }.

Goal. Encode an N-character genome: `ATAGATGCATAG...`

**Standard ASCII encoding.**
- 8 bits per char.
- 8N bits.

**Two-bit encoding encoding.**
- 2 bits per char.
- 2N bits.

| char | hex | binary |
|------|-----|----------|
| A | 41 | 01000001 |
| C | 43 | 01000011 |
| T | 54 | 01010100 |
| G | 47 | 01000111 |

| char | binary |
|------|--------|
| A | 00 |
| C | 01 |
| T | 10 |
| G | 11 |

Amazing but true. Initial genomic databases in 1990s did not use such a code!

Fixed-length code. k-bit code supports alphabet of size $2^k$.

---

## Slide 19

### Genomic code

```
public class Genome {

    public static void compress() {
        Alphabet DNA = new Alphabet("ACTG");
        String s = BinaryStdIn.readString();
        int N = s.length();
        BinaryStdOut.write(N);
        for (int i = 0; i < N; i++) {
            int d = DNA.toIndex(s.charAt(i));
            BinaryStdOut.write(d, 2);
        }
        BinaryStdOut.close();
    }

    public static void expand() {
        Alphabet DNA = new Alphabet("ACTG");
        int N = BinaryStdIn.readInt();
        for (int i = 0; i < N; i++) {
            char c = BinaryStdIn.readChar(2);
            BinaryStdOut.write(DNA.toChar(c));
        }
        BinaryStdOut.close();
    }
}
```

Alphabet data type converts between symbols { A, C, T, G } and integers 0–3.

read genomic string from stdin; write to stdout using 2-bit code

read 2-bit code from stdin; write genomic string to stdout

---

## Slide 20

### Genomic code: test client and sample execution

```
public static void main(String[] args)
{
    if (args[0].equals("-")) compress();
    if (args[0].equals("+")) expand();
}
```

**Tiny test case (264 bits)**

```
% more genomeTiny.txt
ATAGATGCATAGCGCATAGCTAGATGTGCTAGC

java BitsDump 64 < genomeTiny.txt
0100000101010100010000010100011101000001010101000100011101000011
0100000101010100010000010100011101000001010101000100011101000011
0101010001000001010001110100000110101010001000001010001110100000
0101010001000111010101000100011101000001101010100010000010100011
01000011
264 bits

% java Genome - < genomeTiny.txt
??          cannot see bitstream on standard output

% java Genome - < genomeTiny.txt | java BinaryDump 64
0000000000000000000000000001000010010001100101101010010001101110100
1000110110001100101111011011010001101000000
104 bits

% java Genome - < genomeTiny.txt | java HexDump 8
00 00 00 21 23 2d 23 74
8d 8c bb 63 40
104 bits

% java Genome - < genomeTiny.txt | java Genome +
ATAGATGCATAGCGCATAGCTAGATGTGCTAGC          compress-expand cycle
                                           produces original input
```

▸ genomic encoding
▸ **run-length encoding**
▸ Huffman compression
▸ LZW compression

---

## Run-length encoding

Simple type of redundancy in a bitstream. Long runs of repeated bits.

```
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1
```

Representation. Use 4-bit counts to represent alternating runs of 0s and 1s: 15 0s, then 7 1s, then 7 0s, then 11 1s.

```
1 1 1 1 0 1 1 1 0 1 1 1 1 0 0 1 1  ⟵ 16 bits (instead of 40)
  15      7      7       11
```

Q. How many bits to store the counts?
A. We'll use 8.

Q. What to do when run length exceeds maximum count?
A. If run length is longer than 255, we intersperse runs of length 0.

---

## Run-length encoding: Java implementation

```java
public class RunLength
{
    private final static int R = 256;

    public static void compress()
    {  /* see textbook */  }

    public static void expand()
    {
        boolean b = false;
        while (!BinaryStdIn.isEmpty())
        {
            char run = BinaryStdIn.readChar();
            for (int i = 0; i < run; i++)
                BinaryStdOut.write(b);
            b = !b;
        }
        BinaryStdOut.close();
    }

}
```

---

## 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 = 3,000 chars.



A typical bitmap, with run lengths for each row

## Slide 25

▸ genomic encoding
▸ run-length encoding
▸ **Huffman compression**
▸ LZW compression

## Slide 26

Use different number of bits to encode different chars.

Ex. Morse code: • • • — — — • • •

Issue. Ambiguity.
SOS ?
IAMIE ?
EEWNI ?
V7O ?

In practice. Use a medium gap to separate codewords.



S is a prefix of V

## Slide 27

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 char to each codeword.
Ex 3. General prefix-free code.

**Codeword table**

| key | value |
|-----|-------|
| ! | 101 |
| A | 0 |
| B | 1111 |
| C | 110 |
| D | 100 |
| R | 1110 |

**Compressed bitstring**

011111110011001000111111100101 ←— *30 bits*
A  B   RA  CA  DA   B   RA  !

**Codeword table**

| key | value |
|-----|-------|
| ! | 101 |
| A | 11 |
| B | 00 |
| C | 010 |
| D | 100 |
| R | 011 |

**Compressed bitstring**

11000111101011100110001111101 ←— *29 bits*
A B   RA  CA  DA   B   RA  !

## Slide 28

Q. How to represent the prefix-free code?
A. A binary trie!
• Chars in leaves.
• Codeword is path from root to leaf.

**Codeword table**

| key | value |
|-----|-------|
| ! | 101 |
| A | 0 |
| B | 1111 |
| C | 110 |
| D | 100 |
| R | 1110 |

Trie representation



**Compressed bitstring**

011111110011001000111111100101 ←— *30 bits*
A  B   RA  CA  DA   B   RA  !

**Codeword table**

| key | value |
|-----|-------|
| ! | 101 |
| A | 11 |
| B | 00 |
| C | 010 |
| D | 100 |
| R | 011 |

Trie representation



**Compressed bitstring**

11000111101011100110001111101 ←— *29 bits*
A B   RA  CA  DA   B   RA  !

## Prefix-free codes: compression and expansion

**Compression.**
- Method 1: start at leaf; follow path up to the root; print bits in reverse.
- Method 2: create ST of key-value pairs.

**Expansion.**
- Start at root.
- Go left if bit is 0; go right if 1.
- If leaf node, print char and return to root.



Codeword table

| key | value |
|-----|-------|
| ! | 101 |
| A | 0 |
| B | 1111 |
| C | 110 |
| D | 100 |
| R | 1110 |

Trie representation

Compressed bitstring
011111110011001000111111100101 ←— *30 bits*
A  B   RA  CA  DA   B   RA  !

Codeword table

| key | value |
|-----|-------|
| ! | 101 |
| A | 11 |
| B | 00 |
| C | 010 |
| D | 100 |
| R | 011 |

Trie representation

Compressed bitstring
11000111101011100100011111101 ←— *29 bits*
A B  R A  C A  D A B  R A  !

---

## Huffman trie node data type

```
private static class Node implements Comparable<Node>
{
    private char ch;    // Unused for internal nodes.
    private int freq;   // Unused for expand.
    private final Node left, right;

    Node(char ch, int freq, Node left, Node right)
    {
        this.ch    = ch;
        this.freq  = freq;
        this.left  = left;
        this.right = right;
    }

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

    public int compareTo(Node that)
    {   return this.freq - that.freq;  }
}
```

---

## Prefix-free codes: expansion

```
public void expand()
{
    Node root = readTrie();              ←— read in encoding trie
    int N = BinaryStdIn.readInt();       ←— read in number of chars

    for (int i = 0; i < N; i++)
    {
        Node x = root;
        while (!x.isLeaf())              ←— expand codeword for iᵗʰ char
        {
            if (BinaryStdIn.readBoolean())
                x = x.left;
            else
                x = x.right;
        }
        BinaryStdOut.write(x.ch);
    }
    BinaryStdOut.close();
}
```

---

## Prefix-free codes: how to transmit

**Q.** How to write the trie?

**A.** Write preorder traversal of trie; mark leaf and internal nodes with a bit.



*preorder traversal*

leaves

A      D      !      C      R      B

0101000010010101000100010000101010100001101010100101010000010

1        23       4              5   ←— *internal nodes*

Using preorder traversal to encode a trie as a bitstream

```
private static Node readTrie()
{
    if (BinaryStdIn.readBoolean())
    {
        char c = BinaryStdIn.readChar();
        return new Node(c, 0, null, null);
    }
    Node x = readTrie();
    Node y = readTrie();
    return new Node('\0', 0, x, y);
}
```
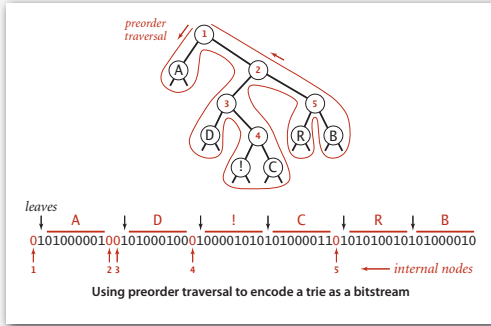
**Note.** If message is long, overhead of transmitting trie is small.

## Prefix-free codes:  how to transmit

Q.  How to read in the trie?

A.  Reconstruct from preorder traversal of trie.



Using preorder traversal to encode a trie as a bitstream

```
private static void writeTrie(Node x)
{
    if (x.isLeaf())
    {
        BinaryStdOut.write(true);
        BinaryStdOut.write(x.ch);
        return;
    }
    BinaryStdOut.write(false);
    writeTrie(x.left);
    writeTrie(x.right);
}
```

## Huffman codes

Q. How to find best prefix-free code?

A.  Huffman algorithm.

David Huffman

Huffman algorithm (to compute optimal prefix-free code):

- Count frequency $freq[i]$ for each char $i$ in input.
- Start with one node corresponding to each char $i$ (with weight $freq[i]$).
- Repeat until single trie formed:
  - select two tries with min weight $freq[i]$ and $freq[j]$
  - merge into single trie with weight $freq[i]$ + $freq[j]$

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

## Constructing a Huffman encoding trie



two tries with
smallest weight

new parent
for those tries

## Constructing a Huffman encoding trie

**Trie representation**



*3 occurrences of w in input*

*1 (LF) 1 (b)*

*labels on path from root are 11010 so 11010 is code for m*

Huffman code for the character stream "it was the best of times it was the worst of times LF"

**Codeword table**

| key | value |
|-----|-------|
| LF  | 101010 |
| SP  | 01 |
| a   | 11011 |
| b   | 101011 |
| e   | 000 |
| f   | 11000 |
| h   | 11001 |
| i   | 1011 |
| m   | 11010 |
| o   | 0011 |
| r   | 10100 |
| s   | 100 |
| t   | 111 |
| w   | 0010 |

```
private static Node buildTrie(int[] freq)
{
    MinPQ<Node> pq = new MinPQ<Node>();
    for (char i = 0; i < R; i++)
        if (freq[i] > 0)
            pq.insert(new Node(i, freq[i], null, null));

    while (pq.size() > 1)
    {
        Node x = pq.delMin();
        Node y = pq.delMin();
        Node parent = new Node('\0', x.freq + y.freq, x, y);
        pq.insert(parent);
    }

    return pq.delMin();
}
```

initialize PQ with singleton tries

merge two smallest tries

not used    total frequency    two subtrees

## Huffman encoding summary

**Proposition.** [Huffman 1950s] Huffman algorithm produces an optimal prefix-free code.

**Pf.** See textbook.

<span style="color:red">no prefix-free code uses fewer bits</span>

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

**Running time.** Using a binary heap ⇒ O(N + R log R).

<span style="color:red">input size</span>   <span style="color:red">alphabet size</span>

**Q.** Can we do better? [stay tuned]

---

Abraham Lempel   Jacob Ziv

---

## 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.

---

## Lempel-Ziv-Welch compression

**LZW compression.**
- Create ST associating W-bit codewords with string keys.
- Initialize ST with codewords for single-character 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 character in the input.

| input | A | B | R | A | C | A | D | A | B | R | A | B | R | A | B | R | A | EOF |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| matches | A | B | R | A | C | A | D | A B | | R A | | B R | | A B R | | | A | |
| output | 41 | 42 | 52 | 41 | 43 | 41 | 44 | 81 | | 83 | | 82 | | 88 | | | 41 | 80 |

**codeword table**

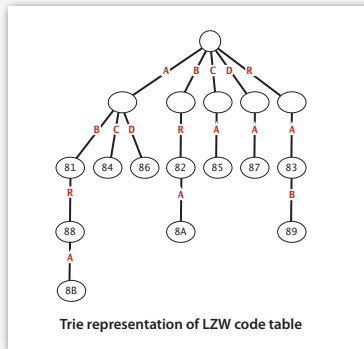| key | value |
|---|---|
| A B | 81 |
| B R | 82 |
| R A | 83 |
| A C | 84 |
| C A | 85 |
| A D | 86 |
| D A | 87 |
| A B R | 88 |
| R A B | 89 |
| B R A | 8A |
| A B R A | 8B |

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

## Representation of LZW code table

Q. How to represent LZW code table?

A. A trie: supports efficient prefix match.

Trie representation of LZW code table

Remark. Every prefix of a key in encoding table is also in encoding table.

## LZW compression: Java implementation

```java
public static void compress()
{
  String input = BinaryStdIn.readString();          // read in input as a string

  TST<Integer> st = new TST<Integer>();             // codewords for single-
  for (char i = 0; i < R; i++)                       // char, radix R keys
      st.put("" + i, (int) i);
  int code = R+1;

  while (input.length() > 0)
  {
      String s = st.prefix(input);                  // find max prefix match s
      BinaryStdOut.write(st.get(s), W);             // write W-bit codeword for s
      int t = s.length();
      if (t < input.length() && code < L)
          st.put(input.substring(0, t+1), code++);  // add new codeword
      input = input.substring(t);                   // scan past s in input
  }

  BinaryStdOut.write(R, W);                          // write last codeword
  BinaryStdOut.close();                             // and close input stream
}
```

## LZW expansion

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

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

## LZW expansion: tricky situation

Q. What to do when next codeword is not yet in ST when needed?

## 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]

## 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

## 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*

## 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.