# Data Compression

## Introduction

### Reasons

- save space when storing
- save time when transmitting

> **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,  ...

### Lossless compression and expansion

- message---binary data $B$ we want to compress
- compress---generate a "compressed" representation $c(B)$
- expand---reconstructs original bitstream $B$
- compression ratio---Bits in $C(B)$/bits in $B$

### Reading and writing binary data

- BinaryStdIn

```
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
```

- BinaryStdOut

```
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
```

# Universal data compression

- No algorithm can compress every bitstream

  - proof by contradiction

    - Suppose you have a universal data compression algorithm $U$ that can compress every bitstream.
    - Given bitstring $B_0$, compress it to get smaller bitstring $B_1$.
    - Compress $B_1$ to get a smaller bitstring $B_2$.
    - Continue until reaching bitstring of size $0$.
    - Implication: all bitstrings can be compressed to $0$ bits!

  - proof by counting

    - Suppose your algorithm that can compress all $1{,}000$-bit strings.
    - $2^{1000}$ possible bitstrings with $1{,}000$ bits.
    - Only $1 + 2 + 4 + \ldots + 2^{998} + 2^{999}$ can be encoded with $\leq 999$ bits.
    - Similarly, only $1$ in $2^{499}$ bitstrings can be encoded with $\leq 500$ bits!

- Undecidability---Impossible to find the best compression algorithm

# Run-length Coding

```
1   public class RunLength {
2       private static final int R    = 256;
3       private static final int LG_R = 8;
4
5       // Do not instantiate.
6       private RunLength() { }
7
8       public static void expand() {
9           boolean b = false;
10          while (!BinaryStdIn.isEmpty()) {
11              int run = BinaryStdIn.readInt(LG_R);
12              for (int i = 0; i < run; i++)
13                  BinaryStdOut.write(b);
14              b = !b;
15          }
16          BinaryStdOut.close();
17      }
18
19      public static void compress() {
20          char run = 0;
21          boolean old = false;
22          while (!BinaryStdIn.isEmpty()) {
23              boolean b = BinaryStdIn.readBoolean();
24              if (b != old) {
25                  BinaryStdOut.write(run, LG_R);
26                  run = 1;
```

```
27                old = !old;
28            }
29            else {
30                if (run == R-1) {
31                    BinaryStdOut.write(run, LG_R);
32                    run = 0;
33                    BinaryStdOut.write(run, LG_R);
34                }
35                run++;
36            }
37        }
38        BinaryStdOut.write(run, LG_R);
39        BinaryStdOut.close();
40    }
41
42    public static void main(String[] args) {
43        if      (args[0].equals("-")) compress();
44        else if (args[0].equals("+")) expand();
45        else throw new IllegalArgumentException("Illegal command line
    argument");
46    }
47
48 }
```

# Huffman Compression

> produce an optimal prefix-free code

## Variable-length Codes

- How to avoid ambiguity?
    - Ensure that no codeword is a prefix of another
        - fixed-length code
        - append special stop char to each codeword
        - general prefix-free code
- How to represent the pre-fix code?
    - A binary trie!
        - Chars in leaves
        - Codewords is path from root to leaf

## Shannon-Fano Codes

> How to find best prefix-free code?

- Shannon-Fano algorithm
    - Partition symbols $S$ into two subsets $S_0$ and $S_1$ of (roughly) equal frequency
    - Codewords for symbols in $S_0$ start with 0; for symbols in $S_1$ start with 1

- Return in $S_0$ and $S_1$

# Huffman Algorithm

- 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 mim weight freq$[i]$ and freq$[j]$
  - merge into single trie with weight freq$[i]$+freq$[j]$

# Implementation

## Huffman Trie

```
1      private static class Node implements Comparable<Node> {
2          private final char ch;
3          private final int freq;
4          private final Node left, right;
5
6          Node(char ch, int freq, Node left, Node right) {
7              this.ch    = ch;
8              this.freq  = freq;
9              this.left  = left;
10             this.right = right;
11         }
12
13         // is the node a leaf node?
14         private boolean isLeaf() {
15             assert ((left == null) && (right == null)) || ((left != null) &&
    (right != null));
16             return (left == null) && (right == null);
17         }
18
19         // compare, based on frequency
20         public int compareTo(Node that) {
21             return this.freq - that.freq;
22         }
23     }
```

## Build Trie

```
1      private static Node buildTrie(int[] freq) {
2
3          // initialize priority queue with singleton trees
4          MinPQ<Node> pq = new MinPQ<Node>();
5          for (char c = 0; c < R; c++)
6              if (freq[c] > 0)
7                  pq.insert(new Node(c, freq[c], null, null));
8
```

```
 9            // merge two smallest trees
10          while (pq.size() > 1) {
11              Node left  = pq.delMin();
12              Node right = pq.delMin();
13              Node parent = new Node('\0', left.freq + right.freq, left, right);
14              pq.insert(parent);
15          }
16          return pq.delMin();
17      }
```

## Expansion

> Running time---linear in input size $N$

```
 1      public static void expand() {
 2
 3          // read in Huffman trie from input stream
 4          Node root = readTrie();
 5
 6          // number of bytes to write
 7          int length = BinaryStdIn.readInt();
 8
 9          // decode using the Huffman trie
10          for (int i = 0; i < length; i++) {
11              Node x = root;
12              while (!x.isLeaf()) {
13                  boolean bit = BinaryStdIn.readBoolean();
14                  if (bit) x = x.right;
15                  else     x = x.left;
16              }
17              BinaryStdOut.write(x.ch, 8);
18          }
19          BinaryStdOut.close();
20      }
```

## Complete Code

```
 1  public class Huffman {
 2
 3      // alphabet size of extended ASCII
 4      private static final int R = 256;
 5
 6      // Do not instantiate.
 7      private Huffman() { }
 8
 9      // Huffman trie node
10      private static class Node implements Comparable<Node> {}
11
12      public static void compress() {
13          // read the input
14          String s = BinaryStdIn.readString();
```

```java
        char[] input = s.toCharArray();

        // tabulate frequency counts
        int[] freq = new int[R];
        for (int i = 0; i < input.length; i++)
            freq[input[i]]++;

        // build Huffman trie
        Node root = buildTrie(freq);

        // build code table
        String[] st = new String[R];
        buildCode(st, root, "");

        // print trie for decoder
        writeTrie(root);

        // print number of bytes in original uncompressed message
        BinaryStdOut.write(input.length);

        // use Huffman code to encode input
        for (int i = 0; i < input.length; i++) {
            String code = st[input[i]];
            for (int j = 0; j < code.length(); j++) {
                if (code.charAt(j) == '0') {
                    BinaryStdOut.write(false);
                }
                else if (code.charAt(j) == '1') {
                    BinaryStdOut.write(true);
                }
                else throw new IllegalStateException("Illegal state");
            }
        }

        // close output stream
        BinaryStdOut.close();
    }

    // build the Huffman trie given frequencies
    private static Node buildTrie(int[] freq) {}

    // write bitstring-encoded trie to standard output
    private static void writeTrie(Node x) {
        if (x.isLeaf()) {
            BinaryStdOut.write(true);
            BinaryStdOut.write(x.ch, 8);
            return;
        }
        BinaryStdOut.write(false);
        writeTrie(x.left);
        writeTrie(x.right);
    }
```

```java
67
68       // make a lookup table from symbols and their encodings
69       private static void buildCode(String[] st, Node x, String s) {
70           if (!x.isLeaf()) {
71               buildCode(st, x.left,  s + '0');
72               buildCode(st, x.right, s + '1');
73           }
74           else {
75               st[x.ch] = s;
76           }
77       }
78
79       public static void expand() {}
80
81       private static Node readTrie() {
82           boolean isLeaf = BinaryStdIn.readBoolean();
83           if (isLeaf) {
84               return new Node(BinaryStdIn.readChar(), -1, null, null);
85           }
86           else {
87               return new Node('\0', -1, readTrie(), readTrie());
88           }
89       }
90
91       public static void main(String[] args) {
92           if      (args[0].equals("-")) compress();
93           else if (args[0].equals("+")) expand();
94           else throw new IllegalArgumentException("Illegal command line
   argument");
95       }
96
97   }
```

# LZW Compression

> Lempel-Ziv-Welch compression

## Steps for 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

## 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

## Compression

```java
public static void compress() {
    String input = BinaryStdIn.readString();
    TST<Integer> st = new TST<Integer>();

    // since TST is not balanced, it'd be better to insert in a different order
    for (int i = 0; i < R; i++)
        st.put("" + (char) i, i);

    int code = R+1;  // R is codeword for EOF

    while (input.length() > 0) {
        String s = st.longestPrefixOf(input);  // Find max prefix match s.
        BinaryStdOut.write(st.get(s), W);       // Print s's encoding.
        int t = s.length();
        if (t < input.length() && code < L)    // Add s to symbol table.
            st.put(input.substring(0, t + 1), code++);
        input = input.substring(t);            // Scan past s in input.
    }
    BinaryStdOut.write(R, W);
    BinaryStdOut.close();
}
```

## Expansion

```
 1      public static void expand() {
 2          String[] st = new String[L];
 3          int i; // next available codeword value
 4
 5          // initialize symbol table with all 1-character strings
 6          for (i = 0; i < R; i++)
 7              st[i] = "" + (char) i;
 8          st[i++] = "";                          // (unused) lookahead for EOF
 9
10          int codeword = BinaryStdIn.readInt(W);
11          if (codeword == R) return;             // expanded message is empty string
12          String val = st[codeword];
13
14          while (true) {
15              BinaryStdOut.write(val);
16              codeword = BinaryStdIn.readInt(W);
17              if (codeword == R) break;
18              String s = st[codeword];
19              if (i == codeword) s = val + val.charAt(0);    // special case hack
20              if (i < L) st[i++] = val + s.charAt(0);
21              val = s;
22          }
23          BinaryStdOut.close();
24      }
```

## Complete Implementation

```
 1  public class LZW {
 2      private static final int R = 256;        // number of input chars
 3      private static final int L = 4096;       // number of codewords = 2^W
 4      private static final int W = 12;         // codeword width
 5
 6      // Do not instantiate.
 7      private LZW() { }
 8
 9      public static void compress() {}
10
11      public static void expand() {}
12
13      public static void main(String[] args) {
14          if      (args[0].equals("-")) compress();
15          else if (args[0].equals("+")) expand();
16          else throw new IllegalArgumentException("Illegal command line argument");
17      }
18  }
```

## Other versions

- LZ77
- LZ78
- Deflate/zlib = LZ77+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 | 7.00 |
| 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 |

## Summary

- Lossless compression
  - Huffman---represent fixed-length symbols with variant-length codes
  - LZW---represent variable-lengh symbols with fixed-length codes
- Lossy compression
  - JPEG, MPEG, MP3,...
  - FFT, wavelets, fractals,...