Correction From Tuesdays' Lecture -Checking the CRC

- Remainder = 0
- Message is Correct
- The quotient is discarded; only the fact that the remainder is zero or non-zero is important

Data Compression

- Introduction
- Huffman
- LZW
- Burrows Wheeler Transform

Compression Techniques

- Variable length codes
 - Huffman, or similar variable length codes
 - Assign short codes to more frequent symbols
 - Needs a coding dictionary
 - Not the most effective
- Run-length coding
 - Many files have lots of repeating characters, for example faxes
 - 0000000000000011111100000000111 may be encoded as 14,0,7,1,8,0,3,1
- Codebooks
 - Assign codes to frequent words and then run as an extension to VL codes
 - Excellent in some contexts but inflexible

Compression Techniques

Dictionary-based

- Use redundancy in text to replace already known letter sequences or *phrases* by dictionary indices
- These compressors learn from the text and adapt to it
- Relative or Differential coding
 - Used for image data
 - Sends the changes between frames
- Statistical-based
 - Assign each symbol a probability
 - Encode the more probable symbol using few bits

Huffman Encoding

- Obtain a symbol frequency count for the message
- Build the Huffman tree
- Encode the message
- If the decoder does not have access to the same code-set, this will need to be transmitted with the message

Huffman Encoding - Example

- Assume a message has the following frequencies
 - {a=4,e=12,i=8,o=2,u=2,y=1}
- Assume every symbol is a single leaf node
- Build a binary tree by continually adding nodes using the two lowest frequencies
- The first node will consist of u,y

Huffman Tree - Example



Huffman Tree

More frequent symbols have shorter codes
The decoder needs to have access to this code set

LZW Compression

- LZW was proposed by Welch as a variant of a compressor originally proposed by Ziv & Lempel, and is the compressor normally used in data communications
- Uses a dictionary of known phrases or letter sequences
 - The dictionary is initialized with all possible character codes, e.g. 0...255
- Compressor accepts input symbols as long as the most recent symbols correspond to known phrases in the dictionary

LZW Compression

- As soon as the matching fails, the compressor emits the dictionary index of the known phrase and then creates a new phrase consisting of the known phrase followed by the character which 'failed'
- Continues building a new test string for phrase comparisons
- The dictionary is initialized to all values, 0...255

- Encoder input is "the_theme_that_they_heard"
- Assume an empty dictionary

the_t	the_theme_that_they_heard Dictionary						
Pass	Input	Test	Emit	Make	Comment	Index	Contents
		string		entry			
1	t	t				0	256
2							ASCII
3						255	codes
4						256	NA
5						257	NA
6						258	NA
7						259	NA
8						260	NA
9						261	NA
10						262	NA

Encoder input is "the_theme_that_they_heard"

the_t	heme_t		Dictionary				
Pass	Input	Test	Emit	Make	Comment	Index	Contents
		string		entry			
1	t	t			t exists	0	256
2	h	th	t	th			ASCII
3	е	he	h	he		255	codes
4	_	e_	е	e_		256	th
5	t	_t	_	_t		257	he
6	h	th			the exists – continue	258	e_
7	е	the	th(256)	the		259	_t
8	m	em	е	em		260	the
9	е	me	m	me		261	em
10	_	e_			e_exists	262	me

the_theme_that_they_heard Dictionary								
Pass	Input	Test	Emit	Make	Comment	Index	contents	
		string		Entry				
11	t	e_t	e_(258)	e_t		263	e_t	
12	h	th			th exists	264	tha	
13	а	tha	th(256)	tha		265	at	
14	t	at	а	at		266	t_	
15	_	t_	t	t_		267	_th	
16	t	_t			_t exists	268	hey	
17	h	_th	_t(259)	_th		269	У_	
18	е	he			he exists			
19	У	hey	he(257)	hey				
20	_	У_	У	У_				

the_t	heme_t	Dictionary					
Pass	Input	Test	Emit	Make	Comment	Index	Contents
		string		Entry			
21	h	_h	_	_h		270	_h
22	e	he			he exists	271	hea
23	а	hea	he(257)	hea		272	ar
24	r	ar	а	ar		273	rd
25	d	rd	r	rd			

LZW Decoding

The basic steps are

- Remember the previous phrase
- Receive the current index and get its current phrase
- Emit the current phrase
- Create a new dictionary entry of the previous phrase followed by the first symbol of the current phrase

the_them	ie_that_they		Die	ctionary	
Pass	Receive codeword	Make entry	deliver	Index	Contents
1	t		t	0	256
2	h	th	h		ASCII
3	е	he	e	255	codes
4		e		256	th
5	(256) th	_t	th	257	he
6	е	the	е	258	e
7	m	em	m	259	t
8	(258) e_	me	e	260	the
9	(256) th	e_t	th	261	em
10	а	tha	а	262	me
11	t	at	t	263	e_t
12	(259) _t	t_	_t	264	tha

Summary of VL and dictionary-based Codes

- Huffman maps fixed length symbols to variable length codes
 - Can be made adaptive
 - Optimal when symbol frequencies are powers of two
- LZW is a dictionary-based compression method
 - Builds a growing dictionary of phrases which are mapped to a sequentially increasing codeset

BZIP2 Compressor (Burrows-Wheeler Transform)

- Write out all cyclic permutations of the input as rows of a matrix
- Sort the rows into alphabetic order, this collects together similar contexts
- Transmit the last symbol of each row, and the index of the first row
 - Similar contexts produce similar symbols and often runs of the one symbol

BZIP2 Compressor (Burrows-Wheeler Transform)

- Use a Move-To-Front algorithm to recode each symbol
 - Replace each by the number of different symbols seen since it was last encountered
 - This produces a skewed symbol distribution

Value	0	1	2	3	4	5	6	7
Frequency	66.8%	9.0%	4.0%	2.9%	2.3%	1.8%	1.6%	1.4%

Use a Huffman encoder on the MTF output

BWT Example

- Get all permutations
- Sort them
- Transmit the last column along with the index of the original string to the MTF

String
takaka
akakat
kakata
akatak
kataka
atakak

Row	String
1	akakat
2	akatak
3	atakak
4	kakata
5	kataka
6	takaka



Reversing BWT

- We have the last column, and therefore can build up the first column of the matrix
- Sorting the pairs reveals the second column
- The transmitted index reveals the original string

Row	String	Row	String
1	a????t	1	ak???t
2	a????k	2	ak???k
3	a????k	3	ak???k
4	k????a	4	ka???a
5	k????a	5	ka???a
6	t????a	6	ta???a

MTF Example

- This output is then transmitted to the Huffman encoder
- Is easily reversed
- Benefits become apparent when using a longer string

 - 0000000010000000000000000

MTF Example

- Initialize the MTF alphabet 0..#A
- After encoding each symbol, 'push' it to the front of the MTF alphabet

Symbol	Code
А	0
К	1
t	2

Symbol	Before			After			Emitted
	MTF			MTF			Code
	0	1	2	0	1	2	
t	а	k	t	t	а	k	2
k	t	а	k	k	t	а	2
k	k	t	а	k	t	а	0
а	k	t	а	а	k	t	2
а	а	k	t	а	k	t	0
а	а	k	t	а	k	t	0

Compression Results

LZW - 4 - 4.5 bits per character or 50 - 55% of the original text size
BZIP 2 - 2.5 bits per character or 25 - 30 % of the original text size

Recommended Reading

 Understanding Data Communications and Networks
 Pages 189 – 194 (Huffman Codes)
 Liv-Zempel Codes (LZW)