

**ĐẠI HỌC THÁI NGUYÊN
TRƯỜNG ĐH CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG**

NGUYỄN THỊ PHƯỢNG

**KỸ THUẬT NÉN DỮ LIỆU
BURROW WHEELER VÀ CÁC CẢI TIẾN**

LUẬN VĂN THẠC SỸ: KHOA HỌC MÁY TÍNH

**ĐẠI HỌC THÁI NGUYÊN
TRƯỜNG ĐH CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG**

NGUYỄN THỊ PHƯỢNG

**KỸ THUẬT NÉN DỮ LIỆU
BURROW WHEELER VÀ CÁC CẢI TIẾN**

**CHUYÊN NGÀNH: KHOA HỌC MÁY TÍNH
MÃ SỐ CHUYÊN NGÀNH: 60 80 01**

LUẬN VĂN THẠC SĨ: KHOA HỌC MÁY TÍNH

**HƯỚNG DẪN KHOA HỌC:
PGS TSKH NGUYỄN XUÂN HUY**

Thái Nguyên - 2011

LỜI CAM ĐOAN

Tôi xin cam đoan: Luận văn “*Kỹ thuật nén dữ liệu Burrow Wheeler và các cải tiến*” là công trình nghiên cứu của riêng tôi. Các số liệu, kết quả nghiên cứu trong luận văn được sử dụng trung thực và có nguồn trích dẫn.

LỜI CẢM ƠN

Em xin gửi lời cảm ơn sâu sắc tới thầy PGS TSKH Nguyễn Xuân Huy - Viện Công nghệ thông tin, người đã gợi mở và định hướng cho em tìm hiểu về lĩnh vực giấu tin trong ảnh. Thầy đã hết lòng giúp đỡ, tạo điều kiện cho em nghiên cứu và hoàn thành luận văn này.

Em xin cảm ơn các thầy cô trong Viện Công nghệ thông tin, các thầy cô giáo khoa Công nghệ thông tin ĐH Thái nguyên, đã giảng dạy và giúp đỡ em trong hai năm học qua.

Cuối cùng tôi xin cảm ơn tới gia đình, các bạn cùng lớp và các bạn đồng nghiệp đã giúp đỡ, động viên, cùng nghiên cứu, đóng góp ý kiến, chia sẻ kinh nghiệm với tôi trong suốt quá trình học tập và làm luận văn!

Thái Nguyên - 2011

Nguyễn Thị Phượng

MỤC LỤC	i
DANH MỤC CÁC KÝ HIỆU, CÁC CHỮ VIẾT TẮT	iii
DANH SÁCH BẢNG BIỂU	iv
DANH SÁCH HÌNH ẢNH	v
MỞ ĐẦU	1
1. Lý do chọn đề tài	1
2. Mục tiêu của đề tài	2
3. Đối tượng và phạm vi nghiên cứu	2
4. Hướng nghiên cứu của đề tài	3
5. Những nội dung nghiên cứu chính	3
6. Phương pháp nghiên cứu	3
7. Ý nghĩa khoa học của đề tài	3
CHƯƠNG 1: TỔNG QUAN VỀ NÉN DỮ LIỆU	4
1.1. Nén dữ liệu	4
1.1.1. Khái niệm về dữ liệu.....	4
1.1.2. Sự trùng lặp dữ liệu.....	4
1.1.3. Nén dữ liệu.....	5
1.2 Các phương pháp nén dữ liệu cơ bản	5
1.2.1. Nén không tổn hao.....	5
1.2.2. Nén tổn hao.....	6
1.3. Dữ liệu ký hiệu và các mã	7
1.3.1. Dữ liệu kí hiệu.....	12
1.3.2. Mã chiều dài thay đổi.....	7
1.3.3. Mã tiền tố và các cây nhị phân.....	8
1.4. Cơ bản về lý thuyết thông tin	11
1.5. Đơn vị đo đặc tính nén	12
CHƯƠNG 2: KỸ THUẬT NÉN DỮ LIỆU BURROWS WHEELER VÀ CÁC CẢI TIẾN	14
2.1. Chuyển đổi Burrows – Wheeler	14
2.1.1. Giới thiệu.....	14
2.1.2. Chuyển đổi Burrows-Wheeler thuận.....	14
2.1.3. Chuyển đổi Burrows-Wheeler nghịch.....	16
2.2. Kỹ thuật nén dữ liệu Burrows-Wheeler	23

2.3. Các cải tiến với kỹ thuật nén dữ liệu Burrows-Wheeler	27
2.3.1. Các định nghĩa	35
2.3.2. Sự đảo ngược tần số (IF)	30
2.3.3. Mã hóa khoảng cách (DC).....	32
2.3.4. Phương pháp đếm trọng số tần số (WFC).....	35
2.3.5. Những thay thế MTF khác	34
2.3.6. Mã hoá Run Length	
2.3.7. Các cải tiến với mã hóa RLE	36
2.3.7.1. Hoạt động chung	36
2.3.7.2. Vị trí mới cho giai đoạn RLE.....	37
2.3.7.3. Thuật toán RLE-EXP	38
2.3.7.4. Thuật toán RLE-BIT	39
2.3.8. Các cải tiến với đảo ngược tần số	41
2.3.8.1. Sắp xếp biểu tượng bằng phân phối tần số	41
2.3.8.2. Thứ tự sắp xếp.....	42
2.3.8.3. Giai đoạn EC	43
2.3.9. Các cải tiến với đếm tần số trọng số	44
2.3.9.1. Phân cấp mịn hơn (Finer Graduation).....	44
2.3.9.2. Tính toán các trọng số.....	44
2.3.9.3. Giai đoạn EC	46
2.3.10. Một thuật toán nén Burrows-Wheeler được cải tiến	47
2.3.10.1. Lựa chọn giai đoạn GST	47
2.3.10.2. So sánh tỉ lệ nén và thời gian nén	49
Kết luận	51
CHƯƠNG 3: CÀI ĐẶT THỬ NGHIỆM	52
3.1. Sơ đồ nén số học kết hợp với BWT và MTF	52
3.1.1. Thuật toán nén	52
3.1.2. Thuật toán giải nén	52
3.2. Cài đặt thử nghiệm	53
3.3. Kết luận	54
KẾT LUẬN VÀ DỰ KIẾN	55
TÀI LIỆU THAM KHẢO	56

DANH MỤC CÁC KÝ HIỆU, CÁC CHỮ VIẾT TẮT

Chữ viết tắt	Diễn giải	Ý nghĩa
AWFC	Advanced Weighted Frequency Count	Đếm trọng số tần số cao cấp
BWCA	Burrows-Wheeler Compression Algorithm	Thuật toán nén Burrows-Wheeler
BWT	Burrows Wheeler Transform	Chuyển đổi Burrows Wheeler
DC	Distance Coding	Mã hóa khoảng cách
EC	Entropy Coding	Mã hóa Entropy
GST	Global Structure Transformation	Chuyển đổi cấu trúc tổng thể
IF	Inversion Frequencies	Sự đảo ngược tần số
IFC	Inversion Frequencies Count	Đếm gia tăng tần số
LUA	List Update Algorithm	Thuật toán cập nhật danh sách
MTF	Move To Front	Di chuyển lên phía trước
RLE	Run Length Encoding	Mã hóa loạt dài
RLE – BIT	Run Length Encoding - BIT	Thuật toán RLE – BIT
RLE – EXP	Run Length Encoding - EXP	Thuật toán RLE – EXPBIT
RLE0	Run Length Encoding 0	Mã hóa chuyển đổi run 0
RMB	RLE Mantissia Buffer	Luồng dữ liệu riêng biệt
SIF	Sort Inversion Frequencies	Sự đảo ngược tần số có sắp xếp
WFC	Weighted Frequency Count	Đếm trọng số tần số

DANH SÁCH BẢNG BIỂU

- Bảng 2.1- 2.2: Mã hóa Move-To-Front
- Bảng 2.3: Sự đảo ngược tần số
- Bảng 2.4.: Mã hóa khoảng cách
- Bảng 2.5: Những giá trị xếp hạng trung bình r_x và thời gian
- Bảng 2.6: Tỷ lệ nén với giai đoạn RLE trước và sau giai đoạn WFC trong bps.
- Bảng 2.7: Các run ngưỡng với $t=2$
- Bảng 2.8: Mã hóa RLE-BIT của chiều dài run
- Bảng 2.9: Tỷ lệ nén theo bps cho các giai đoạn IF và SIF.
- Bảng 2.10: Biểu thị S với mỗi file của Calgary Corpus với $f_a \geq 2F_{avg}$
- Bảng 2.11: Tỷ lệ nén theo bpc với $w_{p_0, p_1, S}(t)$.
- Bảng 2.12: Tỷ lệ nén theo bpc với lược đồ SIF và AWFC.
- Bảng 2.13: Tỷ lệ nén với Calgary Corpus theo bpc
- Bảng 2.14: Thời gian nén và giải nén với Calgary Corpus theo giây

DANH SÁCH HÌNH ẢNH

- Hình 1.1: Minh họa việc trùng lặp dữ liệu của các frame hoạt hình
- Hình 1.2: Máy nén và máy giải nén.
- Hình 1.3: Bộ mã hóa và bộ giải mã
- Hình 1.4: Những thuật toán nén không tổn hao
- Hình 1.5: Những thuật toán nén tổn hao
- Hình 1.6: Mã và dữ liệu nguồn
- Hình 1.7: Đặc tính tiền tố và các cây nhị phân
- Hình 1.8: Không phải mã tiền tố nhưng có thể giải mã duy nhất
- Hình 2.1: Minh họa chuyển đổi BWT thuận
- Hình 2.2: Mảng R được sử dụng để sắp xếp file
- Hình 2.3: Minh họa chuyển đổi BWT nghịch
- Hình 2.4: Sử dụng thứ tự ký tự để thực hiện chuyển đổi ngược
- Hình 2.5: Mảng (A_s) mặc nhiên được khôi phục để giải mã
- Hình 2.6: Các mảng phụ trợ V và W có thể được sử dụng để giải mã xâu mẫu
- Hình 2.7: Lược đồ nén Burrows-Wheeler cơ bản
- Hình 2.8a,b: Minh họa về Mã hóa Huffman
- Hình 2.9 – 2.10: Minh họa mã hóa số học
- Hình 2.11: Lược đồ nén BWT
- Hình 2.12: Sự chia sẻ các Book1
- Hình 2.13: Thuật toán nén Burrows-Wheeler sử dụng giai đoạn RLE0
- Hình 2.14: Minh họa chuyển đổi RLE0
- Hình 2.15: Thuật toán RLE-EXP
- Hình 2.16: Thuật toán RLE-BIT
- Hình 2.17: BWCA với giai đoạn GST hỗn hợp

MỞ ĐẦU

1. Lý do chọn đề tài

Một trong các chức năng chính của máy tính là xử lý và lưu trữ dữ liệu. Bên cạnh việc xử lý nhanh người ta còn quan tâm đến việc lưu trữ được nhiều dữ liệu nhưng lại tiết kiệm được vùng nhớ và giảm chi phí lưu trữ. Về mặt lý thuyết thì các thiết bị lưu trữ là không có giới hạn, nhưng ngày nay do nhu cầu xử lý nhiều tập tin, nhiều loại dữ liệu trong cùng một tệp do vậy mà kích thước tập tin trở nên khá lớn.

Trong nhiều năm gần đây, mạng máy tính đã trở nên khá phổ biến trên thế giới. Sự ra đời của mạng đã thực hiện ước mơ chinh phục khoảng cách của con người. Những lợi ích mà mạng cung cấp rất đa dạng và phong phú trên các lĩnh vực khác nhau của toàn xã hội như cung cấp, trao đổi thông tin giữa các máy tính, giữa máy chủ với server hoặc giữa các server với nhau. Điều này dẫn đến phải làm thế nào để giảm thiểu thời gian, chi phí sử dụng để trao đổi dữ liệu trên mạng. Nó đồng nghĩa với việc bên cạnh nâng cao chất lượng của các thiết bị truyền dữ liệu trên mạng thì mặt khác chúng ta phải nghĩ ra một phương pháp để sao cho việc truyền dữ liệu có hiệu quả hơn.

Tất cả những vấn đề trên nảy sinh ra nhu cầu nén dữ liệu với mong muốn thu gọn kích thước các tập tin làm cho thông tin chiếm không gian trên đĩa là ít nhất. Nó là một trong những kỹ thuật quyết định đến cuộc cách mạng đa phương tiện kỹ thuật số đang diễn ra trong nhiều thập kỷ.

Trong một quá trình phát triển lâu dài, nhiều kỹ thuật nén dữ liệu đã ra đời và được chia làm 2 nhóm kỹ thuật chính là nén tổn hao và nén không tổn hao. Nén không tổn hao là kỹ thuật nén mà sau đó ta có thể khôi phục lại chính xác dữ liệu ban đầu. Nén tổn hao là kỹ thuật nén mà sau khi nén ta không thể khôi phục lại chính xác dữ liệu ban đầu của nó. Nén không tổn hao rất được mong đợi vì tỉ lệ nén cũng như tốc độ nén cao của nó. Tuy nhiên, kỹ thuật này chỉ được dùng cho nén âm thanh, hình ảnh bởi cách thức mà các hệ thống thị giác và thính giác của con người làm việc có thể chấp nhận được. Với dữ liệu gốc của nguồn là rất quan trọng mà ta không thể để mất bất kỳ chi tiết nào dù như hình ảnh y tế, văn bản, các hình ảnh được bảo vệ vì lý do

pháp lý, một số file khả thi của máy tính ... thì ta không thể sử dụng kỹ thuật nén tồn hao.

Nhiều kỹ thuật nén không tồn hao đã ra đời như: Phương pháp mã hóa Entropy bao gồm mã số học và mã Huffman. Sau đó, hàng loạt các kỹ thuật mới ra đời để cải tiến các kỹ thuật trên như: Mã hóa RLE, CD, MTF, LZW... Và gần đây là kỹ thuật nén dữ liệu Burrows Wheeler được công bố bởi Burrows và Wheeler năm 1994. Trong vòng hơn một thập kỷ qua, thuật toán nén Burrows Wheeler [6] đã trở thành một trong những công cụ then chốt trong lĩnh vực nén dữ liệu nói chung. Lý do thành công của nó là tốc độ nén cao kết hợp với tỷ lệ nén tốt. Nhiều cải tiến của kỹ thuật này cũng đã được trình bày. Hiện nay, nén dữ liệu đặc biệt là kỹ thuật nén không tồn hao Burrows Wheeler đang là vấn đề quan tâm rất lớn của các cá nhân, tổ chức, trường học, viện nghiên cứu... trên thế giới.

Chính vì vậy, tôi đã chọn đề tài “*Kỹ thuật nén dữ liệu Burrows Wheeler và các cải tiến*”.

Cấu trúc của luận văn được chia làm 3 chương. **Chương 1: Tổng quan về nén dữ liệu.** Trình bày các khái niệm cơ bản như “dữ liệu”, “Nén dữ liệu”... Các phương pháp nén dữ liệu...**Chương 2: Kỹ thuật nén dữ liệu Burrows Wheeler và các cải tiến.** Trong chương này trình bày cách làm việc của chuyển đổi Burrows Wheeler, kỹ thuật nén dữ liệu Burrows Wheeler và các cải tiến với kỹ thuật nén này. **Chương 3: Cài đặt thử nghiệm.** Áp dụng chuyển đổi BWT tiến hành một kỹ thuật nén số học kết hợp với BWT và MTF trên đối tượng là tệp văn bản. Xây dựng chương trình thử nghiệm áp dụng thuật toán nén số học kết hợp với BWT và MTF.

2. Mục tiêu của đề tài

Luận văn tập trung tìm hiểu kỹ thuật nén BWT và các cải tiến của kỹ thuật. Cuối cùng là cài đặt thử nghiệm chương trình nén số học kết hợp với BWT và MTF.

3. Đối tượng và phạm vi nghiên cứu

Đối tượng và phạm vi nghiên cứu của luận văn tập trung vào kỹ thuật nén dữ liệu Burrows Wheeler và các cải tiến. Từ đó xây dựng chương trình ứng dụng nén số học kết hợp với Burrows Wheeler và MTF, sử dụng phương pháp Quicksort để sắp xếp và áp dụng trên đối tượng dữ liệu là tệp văn bản.