

ĐẠI HỌC THÁI NGUYÊN
TRƯỜNG ĐẠI HỌC CNTT & TRUYỀN THÔNG

Nguyễn Thành Trung

***NÉN DỮ LIỆU KẾT HỢP VỚI CÁC PHƯƠNG PHÁP
BIẾN ĐỔI SƠ BỘ DỮ LIỆU***

Chuyên ngành : KHOA HỌC MÁY TÍNH

Mã số : 60 48 01

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

NGƯỜI HƯỚNG DẪN KHOA HỌC

TS Bùi Văn Thanh

Thái Nguyên, năm 2013

LỜI CAM ĐOAN

Tôi xin cam đoan Luận văn này được thực hiện bằng công sức của mình theo sự hướng dẫn của TS Bùi Văn Thanh, không sao chép từ công trình khác. Mọi thông tin tham khảo đều được trích dẫn đầy đủ. Nếu có gì gian dối tôi xin chịu hoàn toàn trách nhiệm.

Thái Nguyên, tháng 9 năm 2013

Học viên

Nguyễn Thành Trung

LỜI CẢM ƠN

Tôi xin chân thành nói lời cảm ơn Thầy giáo TS Bùi Văn Thanh, người đã tận tình giúp đỡ hướng dẫn tôi trong suốt quá trình thực hiện luận văn cùng với những kinh nghiệm quý báu trong nghiên cứu khoa học cũng như cuộc sống từ Thầy.

Tôi chân thành cảm ơn Trường Đại học Công nghệ Thông tin và Truyền thông - Đại học Thái Nguyên, Khoa Công nghệ Thông tin đã tạo điều kiện tốt nhất cho tôi được học tập và nghiên cứu. Xin cảm ơn quý Thầy giáo, cô giáo đã tận tình giảng dạy, giúp đỡ và hướng dẫn tôi trong suốt khóa học.

Cảm ơn các bạn đồng nghiệp đã động viên đóng góp ý kiến cho luận văn của tôi.

Mặc dù đã cố gắng hết sức cùng với sự tận tâm của thầy giáo hướng dẫn, song do trình độ còn hạn chế nên Luận văn khó tránh khỏi những thiếu sót. Tôi rất mong nhận được sự thông cảm và góp ý của quý thầy cô và các bạn.

Thái Nguyên, tháng 9 năm 2013

Học viên

Nguyễn Thành Trung

MỤC LỤC

	Trang
LỜI CAM ĐOAN	ii
LỜI CẢM ƠN	iii
MỤC LỤC	iv
DANH MỤC CÁC KÝ HIỆU VÀ CÁC CHỮ VIẾT TẮT	vi
DANH MỤC CÁC HÌNH VẼ & BẢNG BIỂU	vii
CHƯƠNG I: TỔNG QUAN VỀ NÉN DỮ LIỆU	1
1.1 Tổng quan về nén dữ liệu	1
1.1.1 Các khái niệm cơ bản.....	1
1.1.1.1 Dữ liệu, thông tin và mã hóa.....	1
1.1.1.2 Cơ bản về lý thuyết thông tin.....	4
1.1.1.3 Sự dư thừa dữ liệu.....	5
1.1.1.4 Nén dữ liệu.....	7
1.1.1.5 Quá trình nén và giải nén	8
1.1.1.6 Tỷ lệ nén	9
1.1.2 Phân loại các phương pháp nén dữ liệu	9
1.1.2.1 Nén có hao hụt	9
1.1.2.2 Nén không hao hụt.....	10
1.1.3 Mô hình và mã hóa	15
1.1.4 Các kết quả cơ bản về nén dữ liệu	18
1.1.5 Tổng quan về các trình nén đang được sử dụng rộng rãi hiện nay	20
1.2 Mã hóa Entropy	21
1.2.1 Mã hóa Huffman	21
1.2.1.1 Quá trình mã hóa.....	25
1.2.1.2 Quá trình giải mã	26
1.2.2 Mã hóa số học	27
1.2.2.1 Mô hình mã hóa số học.....	27
1.2.2.2 Quá trình mã hóa.....	27
1.2.2.3 Quá trình giải mã	29
CHƯƠNG II: NÉN KẾT HỢP VỚI KỸ THUẬT BIẾN ĐỔI DỮ LIỆU	32
2.1 Các kỹ thuật biến đổi dữ liệu cơ bản	32

2.1.1 Kỹ thuật biến đổi Burrow-Wheeler	32
2.1.1.1 Biến đổi BWT thuận	32
2.1.1.2 Biến đổi BWT nghịch	35
2.1.2 Kỹ thuật biến đổi Move-To-Front (MTF).....	38
2.1.2.1 Biến đổi MTF thuận.....	38
2.1.2.2 Biến đổi MTF nghịch.....	41
2.2 Một số cải tiến đối với thuật toán MTF	43
2.2.1 Quá trình mã hóa.....	44
2.2.2 Quá trình giải mã	45
2.3 Mô hình nén kết hợp với BWT&MTF.....	46
2.3.1 Quá trình nén	46
2.3.2 Quá trình giải nén	48
CHƯƠNG III: KẾT QUẢ CÀI ĐẶT THỬ NGHIỆM.....	50
3.1 Dữ liệu mẫu	51
3.2 Kết quả thực nghiệm.....	53
3.2.1 Tỷ lệ nén	53
3.2.2 Thời gian nén và giải nén.....	55
3.3 So sánh và đánh giá kết quả thử nghiệm	57
3.4 Kết luận và hướng phát triển tiếp.....	57
3.4.1 Kết luận:.....	57
3.4.2 Hướng phát triển của đề tài.....	58
TÀI LIỆU THAM KHẢO	59

DANH MỤC CÁC KÝ HIỆU VÀ CÁC CHỮ VIẾT TẮT

TT	Từ viết tắt	Viết đầy đủ
1	ARC	ARithmetic Coding
2	De_ARC	Decoder ARithmetic Coding
3	ASCII	American Standard Code for Information Interchange
4	BIT	BIinary digiT
5	BPS	Bits Per Second
6	BWT	Burrow-Wheeler Transform
7	InvBWT	Invert Burrow-Wheeler Transform
8	EC	Entropy Coding
9	LZW	Lempel-Ziv-Welch
10	MTF	Move-To-Front
11	InvMTF	Invert Move-To-Front
12	RLE	Run Length Encoding

DANH MỤC CÁC HÌNH VẼ & BẢNG BIỂU

Hình 1.1: Mô hình quá trình xử lý thông tin.....	1
Hình 1.2: Dữ liệu trong quá trình mã hóa.....	3
Hình 1.3: Dữ liệu trong quá trình giải mã.....	4
Hình 1.4: Quá trình nén và giải nén.....	8
Hình 1.5: Quá trình truyền file multimedia trên mạng.....	8
Hình 1.6: Nén có hao hụt (lossy compression).....	9
Hình 1.7: Nén không hao hụt (lossless compression).....	10
Bảng 1.1: Quá trình mã hóa từ điển.....	12
Bảng 1.2: Quá trình giải mã từ điển.....	13
Bảng 1.3: Dữ liệu mẫu cho mã hóa Huffman.....	13
Hình 1.8: Cây nhị phân trong mã Huffman ở bảng 1.3.....	14
Hình 1.9: Một dãy các giá trị dữ liệu.....	15
Hình 1.10: Một dãy các giá trị dữ liệu.....	16
Bảng 1.4: Mã với các từ mã chiều dài thay đổi.....	17
Hình 1.11: Minh họa nén theo phương pháp Huffman.....	24
Hình 1.12: Cây Huffman.....	25
Hình 1.13: Mô hình nén theo mã số học.....	27
Hình 1.14: Mô hình quá trình nén.....	28
Bảng 1.5: Quá trình xác định $[b_k, l_k)$ trong mã hóa số học.....	28
Hình 1.15: Mô tả quá trình nén theo bảng 1.5.....	29
Hình 2.1: Quá trình quay chuỗi “LAHABANA”.....	32
Hình 2.2: Kết quả sắp xếp theo thứ tự từ điển.....	33
Hình 2.3: Mô tả quá trình biến đổi BWT nghịch.....	36
Bảng 2.1: Quá trình mã hóa MTF.....	39
Hình 2.4: Quá trình biến đổi MTF thuận chuỗi “HLNBAAAA”.....	39
Bảng 2.2: Quá trình giải mã MTF.....	41
Hình 2.5: Quá trình biến đổi MTF nghịch.....	42
Hình 2.6: Quá trình mã hóa MTF (cải tiến).....	44
Hình 2.7: Quá trình giải mã MTF (cải tiến).....	45
Hình 2.8: Lược đồ nén dữ liệu BWT_MTF_EC.....	46
Hình 2.9: Kết quả thực nghiệm với BWT (quá trình thuận).....	47

Hình 2.10: Kết quả thực nghiệm với MTF (quá trình thuận)	47
Hình 2.11: Lược đồ giải nén dữ liệu BWT_MTF_EC.....	48
Hình 2.12: Kết quả thực nghiệm MTF (quá trình nghịch).....	48
Hình 2.13: Kết quả thực nghiệm BWT (quá trình nghịch).....	49
Hình 3.1: Chương trình cài đặt thử nghiệm	50
Bảng 3.1: Mô tả các tệp tin mẫu để thực nghiệm với The Canterbury Corpus	51
Bảng 3.2: Mô tả các tệp tin mẫu để thực nghiệm với The Large Corpus	52
Bảng 3.3: Mô tả các tệp tin mẫu để thực nghiệm	52
Bảng 3.4: Tỷ lệ nén theo % với Canterbury Corpus.....	53
Bảng 3.5: Tỷ lệ nén theo % với Canterbury Corpus lớn.....	54
Bảng 3.6: Tỷ lệ nén theo %.....	54
Bảng 3.7: Thời gian nén và giải nén theo giây với Canterbury Corpus	55
Bảng 3.8: Thời gian nén và giải nén theo giây với Canterbury Corpus lớn	56
Bảng 3.9: Thời gian nén và giải nén theo giây	56

MỞ ĐẦU

1. Lý do chọn đề tài

Trong thời đại CNTT, nhu cầu trao đổi dữ liệu ngày một tăng, dữ liệu cần chia sẻ, trao đổi có dung lượng ngày một lớn hơn, phức tạp hơn và đa dạng hơn. Việc nén dữ liệu nhằm mục đích làm giảm kích thước của dữ liệu gốc giúp cho việc xử lý dữ liệu nhanh hơn (sao chép, di chuyển, tải lên, tải xuống,...).

Việc nén dữ liệu là tất yếu cần thiết do hai lý do chính sau đây. Thứ nhất là do lưu trữ: dữ liệu sau khi nén có dung lượng nhỏ hơn, do vậy cần ít không gian lưu trữ hơn. Thứ hai là giảm thiểu thời gian truyền: dữ liệu sau khi nén có dung lượng nhỏ hơn nên được truyền nhanh hơn.

Do vậy, cần có các thuật toán nén dữ liệu hiệu quả hơn và nhanh hơn vẫn luôn luôn tăng, nhất là với các ứng dụng trực tuyến (hình ảnh, âm thanh). Thực tế là các thuật toán nén dữ liệu đã được liên tục phát triển kể từ khi giới hạn lý thuyết về hiệu suất nén dữ liệu (được định nghĩa như là tỷ lệ các ký tự mã trên một ký tự nguồn) đã được chứng minh một cách chặt chẽ dựa trên lý thuyết thông tin.

Nhiều kỹ thuật nén (không mất thông tin) đã được phát triển như nhóm các phương pháp với tên gọi 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 (Run Length Encoding), LZW (Lempel-Ziv-Welch),...

Nhưng nhìn chung không có kỹ thuật nén nào có thể áp dụng một cách hiệu quả với tất cả các loại dữ liệu. Có những thuật toán cho hiệu suất nén cao, nhưng cài đặt phức tạp và mất nhiều thời gian nén cũng như giải nén. Chính vì vậy, để đạt được hiệu suất nén cao với thuật toán tương đối đơn giản và thời gian nén/giải nén chấp nhận được cần phối hợp các thuật toán biến đổi sơ bộ dữ liệu vào để chuyển dữ liệu cần nén sang dạng dữ liệu thích hợp, sau đó áp dụng một thuật toán nén phù hợp để tăng hiệu suất nén.

Với ý tưởng trên, chúng ta có thể sử dụng kỹ thuật biến đổi Burrow-Wheeler (BWT) kết hợp với kỹ thuật Move-To-Front (MTF) để xử lý sơ bộ dữ liệu cần nén với mục đích thay đổi tính chất thống kê của dữ liệu sao cho việc áp dụng các thuật

toán nén trở lên dễ dàng hơn, sau đó áp dụng phương pháp nén số học hoặc nén Huffman để được kết quả mong muốn.

2. Nhiệm vụ nghiên cứu

- Tìm hiểu tổng quan về nén dữ liệu, kỹ thuật biến đổi dữ liệu Burrows-Wheeler, kỹ thuật biến đổi dữ liệu Move-To-Front, thuật toán nén số học và thuật toán nén Huffman.

- Xây dựng ứng dụng thử nghiệm phối hợp kỹ thuật biến đổi Burrows-Wheeler và kỹ thuật biến đổi Move-To-Front để sơ chế dữ liệu trước khi sử dụng phương pháp nén số học hoặc nén Huffman và áp dụng trên tất cả các dữ liệu.

3. Đối tượng và phạm vi nghiên cứu

- Tổng quan về nén dữ liệu.
- Tìm hiểu kỹ thuật biến đổi dữ liệu Burrows-Wheeler.
- Tìm hiểu kỹ thuật biến đổi dữ liệu Move-To-Front.
- Tìm hiểu thuật toán nén số học và nén Huffman.
- Cài đặt thử nghiệm chương trình ứng dụng nén số học hoặc nén Huffman kết hợp với Burrows-Wheeler và Move-To-Front.

4. Phương pháp nghiên cứu

Sử dụng các phương pháp nghiên cứu chính sau:

- Phương pháp nghiên cứu lý thuyết.
- Phương pháp thực nghiệm.
- Phương pháp trao đổi khoa học, lấy ý kiến chuyên gia.

5. Ý nghĩa khoa học của đề tài

- Giúp tìm hiểu, đánh giá khái quát về nén dữ liệu. Nén = mô hình + mã hóa.
- Xây dựng được chương trình nén phục vụ cho công tác lưu trữ tại đơn vị đang công tác.