

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

**LÊ VĂN NINH**

**NÉN DỮ LIỆU THEO KỸ THUẬT MOVE – TO - FRONT**

**Chuyên ngành: KHOA HỌC MÁY TÍNH**  
**Mã số: 60 48 01**

**LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH**

**Thái Nguyên - 2011**

## *Lời cảm ơn!*

Em xin chân thành cảm ơn các thầy giáo trong Viện công nghệ thông tin, Viện khoa học và công nghệ Việt Nam; các thầy giáo, cô giáo trường Đại học CNTT & TT - Đại học Thái Nguyên đã tạo điều kiện, giúp đỡ em hoàn thành luận văn này.

Đặc biệt, em xin chân thành cảm ơn PGS.TSKH Nguyễn Xuân Huy, thầy giáo đã giảng dạy và trực tiếp hướng dẫn trong suốt quá trình nghiên cứu và hoàn thành luận văn.

Dù đã có nhiều cố gắng, nhưng chắc chắn luận văn sẽ không tránh khỏi những thiếu sót và hạn chế. Vì vậy, rất mong được sự góp ý, chỉ dẫn của các thầy, cô giáo, bạn bè và đồng nghiệp.

Trân trọng cảm ơn!

*Thái Nguyên, tháng 11 năm 2011*

**Lê Văn Ninh**

<b>DANH MỤC TỪ VIẾT TẮT</b>	
Move – To – Front (Di chuyển lên phía trước)	<i>MTF</i>
Borrows – Wheeler Transform (Chuyển đổi BW)	<i>BWT</i>
Run length (Mã hóa Run Length)	<i>RUL</i>
Inversion Frequencies (Sự đảo ngược tần số)	<i>IF</i>
Distance Coding (Mã hóa khoảng cách)	<i>DC</i>
Weighted Frequency Count (Đếm trọng số tần số)	<i>WFC</i>
Increased Frequency Count (Đếm gia tăng tần số)	<i>IFC</i>
Local to Global Transform (Chuyển đổi từ cục bộ đến tổng thể)	<i>LGT</i>

<b>Hình 1.1:</b>	Máy nén và máy giải nén.
<b>Hình 1.2:</b>	Bộ mã hóa và bộ giải mã
<b>Hình 1.3:</b>	Những thuật toán nén không hao tổn
<b>Hình 1.4:</b>	Các thuật toán nén tổn hao
<b>Hình 1.5:</b>	Dữ liệu về nén
<b>Hình 1.6:</b>	Dữ liệu về giải nén
<b>Hình 1.7:</b>	Dữ liệu ký hiệu về nén
<b>Hình 1.8:</b>	Mã và dữ liệu nguồn
<b>Hình 1.9:</b>	Mã tiền tố
<b>Hình 1.10:</b>	Đặc tính tiền tố và các cây nhị phân
<b>Hình 1.11:</b>	Không phải mã tiền tố nhưng có thể giải mã duy nhất
<b>Hình 1.12:</b>	Các điểm ảnh với các màu giống nhau
<b>Hình 1.13:</b>	Một biểu đồ trong những khoảng xác định
<b>Hình 1.14:</b>	Một số dữ liệu ma trận được tập hợp dọc theo một dòng
<b>Hình 1.15:</b>	Một dãy các frame hoạt hình
<b>Hình 2.1(a)</b>	Mảng $A$ chứa tất cả các phép quay của đầu vào mississippi
<b>Hình 2.1(b)</b>	$A_s$ thu được bằng cách sắp xếp $A$ . Cột cuối của $A_s$ (ký hiệu $L$ ) là đầu ra của BWT
<b>Hình 2.2</b>	Mảng $R$ được sử dụng để sắp xếp file mẫu mississippi
<b>Hình 2.3</b>	Mảng $A_s$ với mississippi. $F$ và $L$ là các cột đầu và cuối tương ứng
<b>Hình 2.4</b>	Sử dụng thứ tự ký tự để thực hiện chuyển đổi ngược
<b>Hình 2.5</b>	Mảng ( $A_s$ ) mặc nhiên được khôi phục để giải mã chuỗi pssmipissii
<b>Hình 2.6</b>	Các mảng phụ trợ $V$ và $W$ có thể được sử dụng để giải mã chuỗi mẫu
<b>Hình 2.7</b>	Một số văn bản được chuyển đổi sinh ra từ “Hamlet” của Shakespeare.
<b>Hình 2.8:</b>	Mã hóa Huffman
<b>Hình 2.9:</b>	Mã hóa Huffman ngược
<b>Hình 2.10:</b>	Xác suất và khoảng con khởi tạo của biểu tượng
<b>Hình 2.11:</b>	Mã hóa số học

MỤC LỤC	PAGE
Lời cảm ơn!	
Danh mục các ký hiệu, viết tắt	
Danh mục các hình vẽ	
Mục lục	
<b>MỞ ĐẦU</b>	<b>3</b>
<b>Chương I. TỔNG QUAN VỀ NÉN DỮ LIỆU</b>	<b>5</b>
<b>1.1. Giới thiệu</b>	<b>5</b>
<i>1.1.1. Một số vấn đề về Nén dữ liệu</i>	<i>6</i>
<i>1.1.2 Nén không tổn hao và nén tổn hao</i>	<i>9</i>
<i>1.1.2.1 Nén không tổn hao</i>	<i>9</i>
<i>1.1.2.2 Nén tổn hao</i>	<i>9</i>
<i>1.1.3 Đơn vị đo đặc tính nén</i>	<i>11</i>
<b>1.2. Mã hóa dữ liệu ký hiệu</b>	<b>13</b>
<i>1.2.1. Thông tin, dữ liệu và các mã</i>	<i>14</i>
<i>1.2.2 Dữ liệu ký hiệu</i>	<i>17</i>
<i>1.2.3 Mã chiều dài thay đổi</i>	<i>19</i>
<i>1.2.4 Cơ bản về lý thuyết thông tin</i>	<i>26</i>
<i>1.2.5 Sự dư thừa</i>	<i>27</i>
<b>Chương II. KỸ THUẬT NÉN DỮ LIỆU BURROWS WHEELER</b>	<b>31</b>
<b>2.1 Chuyển đổi Burrows-Wheeler (BWT)</b>	<b>31</b>

2.1.1 Cách làm việc của chuyển đổi Burrows-Wheeler	33
2.1.1.1 Chuyển đổi Burrows – Wheeler thuận	33
2.1.1.2 Chuyển đổi Burrows – Wheeler nghịch	35
<b>2.1.2 Các mã với chuyển đổi Burrows-Wheeler</b>	<b>41</b>
2.1.2.1 Mã hóa Entropy	42
2.1.2.2 Mã hóa Huffman	42
2.1.2.3 Mã hóa số học	43
2.1.2.4 Mã hóa khoảng cách	46
2.1.2.5 Mã hóa run length	47
2.1.2.6 Các phương pháp đếm tần số	48
<b>2.2 Mã hóa Move – To – Front</b>	<b>49</b>
2.2.1 Mã hóa MTF với các biểu tượng là tập hợp các số nguyên	52
2.2.2 Hiệu suất mã hóa MTF	54
<b>Chương III. GIẢI THUẬT MOVE – TO – FRONT VÀ DEMO</b>	<b>60</b>
<b>3.1. Thuật toán nén dữ liệu Move – To - Front</b>	<b>60</b>
3.1.1. Thuật toán mã hóa	61
3.1.2. Thuật toán giải mã	62
<b>3.2. Thực hiện giải thuật bằng ngôn ngữ C</b>	<b>62</b>
<b>KẾT LUẬN</b>	<b>72</b>
<b>TÀI LIỆU THAM KHẢO</b>	<b>73</b>

## MỞ ĐẦU

### 1. Lý do chọn đề tài

Trong các lĩnh vực của công nghệ thông tin và viễn thông hiện nay, việc truyền tải tin tức là công việc xảy ra thường xuyên. Tuy nhiên, thông tin được truyền tải đi thường rất lớn, điều này gây khó khăn cho công việc truyền tải như: gây tốn kém tài nguyên mạng, tiêu phí khả năng của hệ thống,... Để giải quyết vấn đề đó, các thuật toán nén dữ liệu đã được ra đời.

Các kỹ thuật nén được nhúng ngày càng nhiều trong phần mềm và đã trở thành yêu cầu chung cho hầu hết phần mềm ứng dụng như một lĩnh vực nghiên cứu quan trọng và tích cực trong khoa học máy tính.

Trong kỹ thuật truyền tin nối tiếp, do các bit dữ liệu được truyền đi nối tiếp, lại bị giới hạn về dải thông của kênh truyền và giới hạn về các chuẩn ghép nối... nên tốc độ truyền tin tương đối chậm. Nén dữ liệu trước khi truyền đi là một trong các phương pháp nhằm tăng tốc độ truyền dữ liệu. Nguyên tắc của nén dữ liệu là quá trình mã hóa thông tin dùng ít bit hơn so với thông tin chưa được mã hóa bằng cách dùng một hoặc kết hợp các phương pháp nào đó. Dựa theo nguyên tắc này giúp tránh các hiện tượng kênh truyền bị quá tải và việc truyền tin trở nên kinh tế hơn.

Mặc dù các chương trình nén dữ liệu thường sử dụng kết hợp nhiều thuật toán có độ phức tạp khác nhau nhằm đạt được hiệu quả cao nhất cho dữ liệu được nén để đáp ứng yêu cầu đặt ra. Nhưng nhìn chung không thể có phương pháp nén tổng quát nào cho kết quả tốt đối với tất cả các loại tập tin. Kỹ thuật nén tập tin thường được áp dụng cho các tập tin văn bản (Trong đó có một số kí tự nào đó có xác suất xuất hiện nhiều hơn các kí tự khác), các tập tin ảnh bitmap (Mà có thể có những mảng đồng nhất), các tập tin dùng để biểu diễn âm thanh dưới dạng số hoá và các tín hiệu tương tự khác (các tín hiệu này có thể có các mẫu được lặp lại nhiều lần). Đối với các tập tin nhị phân như tập tin chương trình thì sau khi nén cũng không tiết kiệm được

nhieu. Ngoài ra, trong một số trường hợp để nâng cao hệ số nén người ta có thể bỏ bớt một số thông tin của tập tin (Ví dụ như kỹ thuật nén ảnh JPEG).

Nén dữ liệu theo kỹ thuật Move-To-Front (MTF) là một trong những kỹ thuật nén dữ liệu được thiết kế để cải tiến hiệu quả của kỹ thuật nén mã hóa *entropy*. Nó được sử dụng sau kỹ thuật chuyển đổi Burrows-Wheeler để xếp hạng các biểu tượng theo tần số tương quan của chúng. Mục đích là để đạt được một hiệu suất nén tốt hơn cho mã hóa entropy.

Xuất phát từ ý tưởng đó, tôi đã lựa chọn đề tài “Nén dữ liệu theo kỹ thuật Move – To – Front”.

## **2. Đối tượng nghiên cứu:**

- Kỹ thuật nén dữ liệu Move-To-Front

## **3. Phạm vi nghiên cứu:**

- Tìm hiểu tổng quan về nén dữ liệu
- Nén dữ liệu Burrows-Wheeler
- Nén dữ liệu theo kỹ thuật Move – to – front

## **4. Mục tiêu nghiên cứu**

Luận văn tập trung nghiên cứu, đánh giá về nén dữ liệu theo kỹ thuật Move-To-Front. Vận dụng nén dữ liệu trong một số lĩnh vực đặc thù.

## **5. Ý nghĩa khoa học của đề tài**

- Giúp tìm hiểu, đánh giá khái quát về nén dữ liệu theo kỹ thuật MTF.
- Vận dụng được phương pháp nén dữ liệu theo kỹ thuật MTF trong một số lĩnh vực đặc thù.

## **6. 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 thống kê
- Phương pháp trao đổi khoa học, lấy ý kiến chuyên gia.



## **Chương I. TỔNG QUAN VỀ NÉN DỮ LIỆU**

### **1.1. Giới thiệu**

*Nén dữ liệu* trong ngữ cảnh khoa học máy tính là khoa học để biểu diễn thông tin dưới dạng thu gọn. Nói cách khác nén dữ liệu là việc thực hiện thu gọn kích thước các tập tin hoặc làm cho thông tin lưu trữ chiếm không gian lưu trữ ít nhất. Có nhiều cách để thực hiện điều này tùy vào từng đối tượng cụ thể.

Nén dữ liệu đã trở thành yêu cầu chung cho hầu hết phần mềm ứng dụng như một lĩnh vực nghiên cứu quan trọng và tích cực trong khoa học máy tính. Nếu không có các kỹ thuật nén, Internet sẽ không bao giờ phát triển, TV kỹ thuật số, các kỹ thuật truyền thông di động hoặc truyền thông video đã được phát triển trên thực tế.

Các lĩnh vực ứng dụng có liên quan và được thúc đẩy bởi nén dữ liệu gồm có:

- Các hệ thống truyền thông cá nhân: Fax, thư thoại (voice mail) và điện thoại.
- Các hệ thống máy tính: Cấu trúc bộ nhớ, đĩa và băng.
- Tính toán di động.
- Các hệ thống máy tính phân tán.
- Mạng máy tính, đặc biệt là Internet.
- Sự phát triển đa phương tiện, hình ảnh, xử lý tín hiệu.
- Lưu trữ hình ảnh và hội nghị truyền hình.
- Ti vi kỹ thuật số và truyền hình vệ tinh.
- ...

Nhiều vấn đề trên thực tế đã thúc đẩy nhiều nghiên cứu khác nhau về nén dữ liệu. Tương tự, nghiên cứu về nén dữ liệu cũng đã dựa trên hay được

kích thích bởi các lĩnh vực mới khác. Một phần do phạm vi ứng dụng rộng rãi của nó, nén dữ liệu bao trùm nhiều ngành khoa học và có thể được tìm thấy trong nhiều lĩnh vực khác nhau như: Lý thuyết thông tin; Lý thuyết mã hóa; Mạng máy tính và viễn thông; Xử lý tín hiệu kỹ thuật số; Xử lý ảnh; Đa phương tiện; Bảo mật máy tính...

Trong nén dữ liệu, từ *dữ liệu* có nghĩa là thông tin ở dạng kỹ thuật số mà những chương trình máy tính hoạt động và *nén*, có nghĩa là quá trình loại bỏ dư thừa trong dữ liệu. Cụm từ “*nén dữ liệu*”, có nghĩa là đưa ra các kỹ thuật hay cụ thể hơn là thiết kế những thuật toán hiệu quả nhằm để:

- Biểu diễn dữ liệu theo dạng mà chứa ít dư thừa.
- Loại bỏ dư thừa trong dữ liệu.
- Cài đặt thuật toán nén và giải nén.

### 1.1.1. Một số vấn đề về Nén dữ liệu

Một vấn đề nén liên quan đến việc tìm một thuật toán hiệu quả để loại bỏ dư thừa khác nhau từ một kiểu dữ liệu nhất định. Ví dụ cho một xâu  $s$ , câu hỏi là *dãy các biểu tượng có thể thay thế mà chiếm ít không gian lưu trữ là dãy nào? Giải pháp* cho vấn đề nén là thuật toán nén nhằm đưa ra dãy các biểu tượng chứa ít số lượng bit hơn, cộng với các thuật toán giải nén để phục hồi xâu gốc.

Vậy số lượng bit ít hơn là bao nhiêu? Điều đó phụ thuộc vào những thuật toán nhưng nó cũng phụ thuộc vào sự dư thừa có thể chiết ra từ dữ liệu gốc là bao nhiêu. Dữ liệu khác nhau có thể yêu cầu những kỹ thuật khác nhau để xác định dư thừa và loại bỏ dư thừa trong dữ liệu.

Không có giải pháp nào phù hợp cho tất cả vấn đề nén dữ liệu. Theo các nghiên cứu về nén dữ liệu, ta chủ yếu phải phân tích những đặc tính của dữ liệu đã được nén và hy vọng đưa ra một số mô hình để đạt được sự biểu