

I H C THÁI NGUYÊN
TR NG I H C CNTT & TRUY N THÔNG

PH M THU H NG

NÉN V N B N TI NG VI T THEO HUFFMAN

LU N V N TH C S KHOA H C MÁY TÍNH

Thái Nguyên - 2013

I H C THÁI NGUYÊN
TR NG I H C CNTT & TRUY N THÔNG

PH M THU H NG

NÉN V N B N TI NG VI T THEO HUFFMAN

Chuyên ngành: Khoa h c máy tính

Mã s : 60 48 01

LU N V N TH C S KHOA H C MÁY TÍNH

Ng i h ng d n khoa h c: PGS.TS Nguy n H u i n

Thái Nguyên - 2013

L I C M N

tôi rất vui mừng ngày hôm nay là một sự kiện quan trọng, niềm vui của tôi và gia đình, cũng như sự giúp đỡ nhiệt tình của quý thầy cô, bạn bè tôi hoàn thành luận văn này.

Tôi xin trân trọng cảm ơn:

- PGS.TS Nguyễn Hữu Tiến – Giám đốc Trung tâm tính toán kỹ thuật và công nghệ cao Trường Đại học Khoa học Tự nhiên Hà Nội.
- Các thầy cô trong Hội đồng phản biện.

Cùng với tôi xin chân thành cảm ơn các thầy, cô, các bạn đồng nghiệp và gia đình tôi trong thời gian làm luận văn.

Xin trân trọng cảm ơn quý thầy cô, các bạn!

DANH MỤC CÁC HÌNH

Hình 1. Quy trình nén dữ liệu	3
Hình 2. Xây dựng cây nh phân t b ng mã không d ng ti n t	8
Hình 3. S p x p danh sách các ký t	20
Hình 4. Xây dựng cây Huffman	22
Hình 5. Cây Huffman i n y thành ph n	22
Hình 6. M t tr ng h p xây d ng khác	23
Hình 7. L u gi i mã	24
Hình 8. Ý t ng xây d ng cây theo ph ng pháp Shannon – Fano	26
Hình 9. Xây dựng cây theo ph ng pháp Shannon-Fano	27
Hình 10. Mã hóa b ng ph ng pháp Huffman ng	31
Hình 11. Gi i mã b ng ph ng pháp Huffman ng	33
Hình 12. Quá trình th c hi n nén b ng LZ	43
Hình 13. S nén LZ78	47
Hình 14. S gi i nén LZ78	48
Hình 15. S nén LZW	51
Hình 16. S gi i nén LZW	54
Hình 17. Ph ng pháp MTF (t t)	57
Hình 18. Ph ng pháp MTF (x u).....	57
Hình 19. Ph ng pháp BW tìm chu i sau mã hóa.....	59
Hình 20. Hai cách tìm chu i g c.....	60
Hình 21. Giao di n ch ng trình	62

M C L C

L I C M N	iii
DANH M C CÁC HÌNH.....	iv
M U	1
1. <i>t v n</i>	1
2. <i>i t ng và ph m vi nghiên c u</i>	1
2.1. <i>i t ng</i>	1
2.2. <i>Ph m vi</i>	2
3. <i>H ng nghiên c u c a tài</i>	2
4. <i>Ph ng pháp nghiên c u</i>	2
5. <i>Ý ngh a khoa h c c a lu n v n</i>	2
CH NG 1: T NG QUAN V CÔNG NGH NÉN D LI U.....	3
1.1. <i>S l c v nén d li u</i>	3
1.1.1. <i>Khái ni m v nén d li u</i>	3
1.1.2. <i>Nh ng v n ph i gi i quy t trong nén d li u</i>	4
1.1.3. <i>Phân lo i ch ng trình nén</i>	5
1.1.4. <i>ánh giá ch t l ng c a ch ng trình nén</i>	6
1.2. <i>Mã nén d li u</i>	7
1.2.1. <i>nh ngh a mã hoá</i>	7
1.2.2. <i>Các khái ni m v ký t mã hóa</i>	8
1.2.3. <i>Mã t ng và mã phân tách</i>	13
1.2.4. <i>nh lý mã nén</i>	18
CH NG 2. M T S MÃ NÉN C B N.....	21
2.1. <i>Mã hóa Huffman (Huffman coding)</i>	21
2.1.1. <i>Ph ng pháp mã hóa</i>	21
2.1.2. <i>Thu t toán t o mã Huffman</i>	21
2.1.3. <i>Gi i mã thu t toán Huffman :</i>	25
2.2. <i>Mã hóa Huffman ng (Adaptive Huffman coding)</i>	31
2.2.1. <i>Ph ng pháp mã hóa:</i>	31
2.2.2. <i>Thu t toán nén</i>	31
2.2.3. <i>Thu t toán gi i nén</i>	33
2.3. <i>Thu t toán x lý s l p l i c a xâu (RLE)</i>	36

2.3.1. Ph ng pháp:	36
2.3.2. Thu t toán t o mã	36
2.3.3. Quá trình gi i mã	36
2.4. Mã hóa ki u t i n (<i>Dictionary-based compression</i>)	39
2.4.1. Nguyên lý LZ	39
2.4.2. T i n	40
2.4.3. Quá trình th c hi n khi nén b ng mã LZ	41
2.4.4. Các thu t toán nén LZ	42
2.5. M t s ph ng pháp bi n i (<i>transform</i>)	54
2.5.1. Ph ng pháp y v phía tr c (<i>Move to front</i>):	54
2.5.2. Ph ng pháp Burrows – Wheeler (BW):	56
CH NG 3. XÂY D NG CH NG TRÌNH NÉN TI NG VI T S D NG PH NG PHÁP MÃ HÓA HUFFMAN	59
3.1. B g Ti ng vi t	59
3.2. Quy c bi u di n ký t ti ng Vi t.	59
3.3. Chu n d u Ti ng vi t	60
3.3.1. Unicode	60
3.3.2. TCVN3	60
3.3.3. VNI	60
3.4. Ph ng pháp mã hóa Huffman	60
3.5. Gi i thi u ch ng trình	61
3.5.1. H ng d n s d ng	62
3.5.2. K t qu ki m th ch ng trình	64
K T LU N	65
TÀI LI U THAM KH O	66
PH L C	67

M U

1. t v n

M t trong nh ng ch c n ng chính c a máy tính là x lý d li u và l u tr . Bên c nh vi c x lý nhanh, ng i ta còn quan tâm n vi c l u tr c nhi u d li u nh ng l i ti t ki m c vùng nh và gi m chi phí l u tr . V m t lý thuy t thì các thi t b l u tr là không có gi i h n nh ng ngày nay do nhu c u x lý nhi u t p tin, nhi u lo i d li u trong cùng m t t p do v y mà kích th c t p tr nên khá l n. Nh ng v n trên n y sinh ra khái ni m nén d li u, nén d li u là quá trình làm gi m l ng thông tin “d th a” trong d li u g c, do v y l ng thông tin thu c sau nén th ng nh h n d li u g c r t nhi u. Nén d li u là gi i pháp h p lý nh t nh m m c ích gi m chi phí cho ng i s d ng.

Nh chúng ta c ng ã bi t ti ng Vi t là m t ngôn thu c h th ng ch cái Latin s d ng nhi u d u i kèm v i nguyên âm. V i b ng mã ASCII 8 bit s d ng ph bi n trên máy tính, chúng ta có th mã hóa 256 ký t . a ti ng Vi t vào máy tính, các ph n m m ti ng Vi t hi n nay s d ng m t trong hai ph ng pháp mã hóa : mã d ng s n ho c mã t h p xây d ng trang mã ký t ti ng Vi t. B ng mã ph bi n nh t chúng ta th ng s d ng là b ng mã Unicode th hi n ti ng Vi t. Nh ng b ng mã Unicode yêu c u 16 bit th hi n m t ký t i u này s d n n s lãng phí và d th a d li u.

Vì v y, “ Nén v n b n ti ng Vi t theo Huffman ” c em ch n làm lu n v n t t nghi p c a mình.

2. i t ng và ph m vi nghiên c u

2.1. i t ng

- Các ph n m m nén d li u;
- Các thu t toán nén d li u;
- Các ph ng pháp mã hóa ti ng Vi t;
- H th ng ph n m m nén d li u t ó ng d ng vào nén d li u cho ti ng Vi t.

2.2. Phạm vi

- Các khái niệm cơ bản ký tự mã hóa, các thuật toán nén văn bản. Kiến trúc, chức năng và các thành phần của nén dữ liệu cho bài toán nén văn bản. Tìm hiểu sơ lược về phương pháp mã hóa Huffman.
- Các chức năng chính và quy trình thực thi của bài toán nén dữ liệu;
- Hệ thống chương trình cho bài toán nén dữ liệu;

Vì thời gian có hạn, trong khuôn khổ môn luận văn tốt nghiệp cao học, vì các giới hạn quy định bài toán nén dữ liệu chỉ giới hạn một vài thuật toán nén cơ bản.

3. Hướng nghiên cứu của tài

- Tìm hiểu tổng quan về nén dữ liệu và nghiên cứu một số thuật toán nén cơ bản
- Tìm hiểu bài toán nén dữ liệu, tiến hành phân tích;
- Thu thập các số liệu có liên quan;
- Phân tích, đánh giá thông qua các số liệu thu thập được;
- Cài đặt thực nghiệm.

4. Phương pháp nghiên cứu

- Nghiên cứu các tài liệu và viết tổng quan;
- Phương pháp khảo sát thực tế;
- Phương pháp phân tích và đánh giá các thuật toán;
- Nghiên cứu triển khai thuật toán và thực nghiệm hệ thống.

5. Ý nghĩa khoa học của luận văn

- Bản thân hiểu sâu hơn và áp dụng được các thuật toán nén dữ liệu vào thực tế;
- Triển khai một số thuật toán nén dữ liệu qua đó ứng dụng chính là phương pháp mã hóa Huffman vào tìm hiểu Vi tính;
- Xây dựng các chương trình nén dữ liệu dành cho tìm hiểu Vi tính trên máy tính.

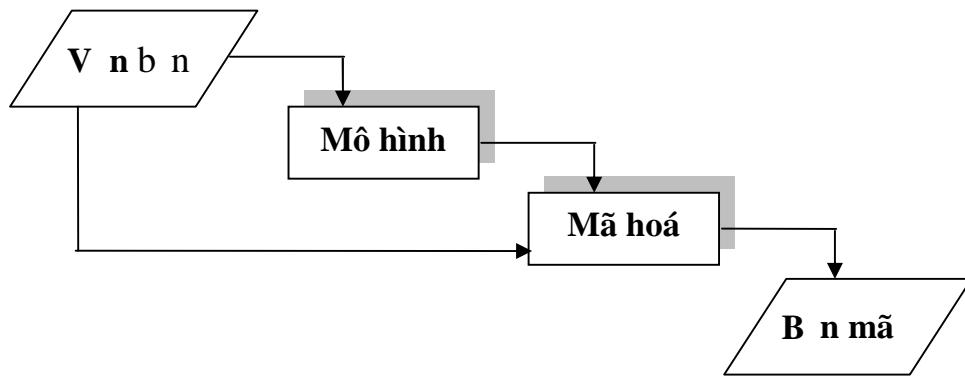
CHƯƠNG 1: TỔNG QUAN VỀ CÔNG NGHỆ NÉN DỮ LIỆU

1.1. Sơ lược về nén dữ liệu

1.1.1. Khái niệm về nén dữ liệu

Nén là một quá trình giảm kích thước không gian cần thiết để biểu diễn cùng một lượng thông tin cho trước. Ngoài ra còn gọi nén là biến đổi một lượng ký hiệu thành một lượng các bit mã.

Quá trình nén như sau:



Hình 1. Quy trình nén dữ liệu

Trong đó:

- Vấn đề là vấn đề ban đầu cần nén.
- Mô hình là tập hợp các cách thức cùng quy tắc sử dụng để xử lý các cách thức vào và đưa ra các bit mã. Một mô hình sẽ xác định chính xác xác suất xuất hiện của từng cách thức và một bit mã sẽ đưa ra các bit mã dựa trên xác suất đó.
- Mã hoá là quá trình thay thế các cách thức trong vấn đề ban đầu bằng các bit mã tương ứng để đưa ra bit mã chính xác.

Như vậy, quá trình nén diễn ra như sau: quá trình mô hình chuyển vào vấn đề cần nén sẽ đưa ra các bit mã. Sau đó, bit mã sẽ được mã hóa và vấn đề ban đầu qua quá trình nén sẽ đưa ra bit mã.

Mã hoá và mô hình là hai giai đoạn hoàn toàn khác nhau vì trong giai đoạn mô hình có rất nhiều cách xử lý các thành phần của văn bản mà cùng số dung lượng thông tin nhưng phương pháp xây dựng mã cho ra các mã.

Nếu văn bản mã có kích thước nhỏ hơn văn bản thì phương pháp nén sẽ có hiệu quả.

Ví dụ :

Chúng ta sẽ dùng cùng phương pháp mã Huffman để mã hoá hai mô hình khác nhau:

- Mô hình 1: dựa trên xác suất của các ký tự để xây dựng mã.
- Mô hình 2: tính toán xác suất phụ thuộc dựa trên ngữ cảnh của ký tự trước đó trong văn bản.

Do mô hình khác nhau nên cùng số dung lượng mã Huffman sẽ mã hoá ra các mã khác nhau.

Tuy nhiên, chúng ta vẫn quen dùng thuật toán mã hoá để cho các quá trình nén văn bản mà dù đó chỉ là một giai đoạn của một quá trình nén.

Ngày nay thuật toán mã hoá thông qua các thuật toán biến đổi ngữ cảnh nào đó.

Có thể có nhiều thuật toán nén dữ liệu khác nhau. Một thuật toán có một kỹ thuật dữ liệu nhất định và cùng một số modem có các kỹ thuật nén thích hợp có nghĩa là chúng có khả năng chọn thuật toán nén thích hợp phụ thuộc vào kỹ thuật dữ liệu cần nén. Trong số các cách mã thì cách nào mã ngắn hơn chúng ta nói là nó nén tốt hơn (so với cách mã khác).

1.1.2. Những vấn đề kỹ thuật trong nén dữ liệu

Một tiêu chuẩn nén dữ liệu là dựa vào thuật toán để giảm thiểu sai lệch phí tổn thông tin khi gửi nhau khi bị gửi đi dữ liệu. Thông thường một chương trình nén cần quan tâm đến khả năng nó có thể xử lý các dữ liệu hay ít dung lượng cần dữ liệu sau khi nén, dữ liệu này còn phụ thuộc vào thuật toán, và dữ liệu gửi vào cần dữ liệu. Những vấn đề dữ liệu gửi vào khác nhau có thể đòi hỏi những thuật toán khác nhau như những vấn đề kỹ thuật chi tiêu sai lệch phí tổn. Cần nhớ rằng