

Mục lục

| | |
|---|-----------|
| Mục lục | i |
| Danh mục các hình ảnh | iv |
| MỞ ĐẦU | 1 |
| I. Đặt vấn đề | 1 |
| II. Nội dung nghiên cứu | 2 |
| III. Bố cục của luận văn | 4 |
| Chương I. TỔNG QUAN VỀ NHẬN DẠNG CHỮ VIẾT | 5 |
| VÀ PHÂN TÍCH TRANG TÀI LIỆU | 5 |
| I.1. Ảnh tài liệu và nhận dạng ảnh tài liệu | 5 |
| I.1.1. Tổng quan về ảnh tài liệu | 5 |
| I.1.2. Nhận dạng tài liệu và vai trò của phân tích ảnh tài liệu | 6 |
| I.2. Cấu trúc của ảnh tài liệu | 7 |
| I.2.1. Cấu trúc vật lý | 8 |
| I.2.2. Cấu trúc logic | 10 |
| I.3. Quá trình phân tích tài liệu | 10 |
| I.3.1. Tiền xử lý(preprocessing): | 11 |
| I.3.2. Phân tích cấu trúc vật lý | 12 |
| I.3.3. Phân tích cấu trúc logic | 13 |
| I.4. Một số hệ thống phân tích tài liệu hiện nay | 14 |
| I.4.1. VnDOCR | 14 |
| I.4.2. OminiPage | 18 |
| I.4.3. Finereader | 20 |
| I.5. Kết luận | 22 |
| Chương II: CÁC PHƯƠNG PHÁP TIẾP CẬN | 23 |
| ĐỀ PHÂN TÍCH TRANG TÀI LIỆU | 23 |
| II.1. Các phương pháp phân tích định dạng trang tài liệu | 23 |
| II.1.1. Top-down | 23 |

| | |
|--|----|
| II.1.2. Bottom-up | 30 |
| II.1.3. Phương pháp Tách và Nối thích nghi (Adaptive Split – and – Merge) . | 32 |
| II.1.4. Fractal Signature (FS)..... | 34 |
| II.2. Lựa chọn giải pháp..... | 38 |
| II.3. Thiết kế hệ thống..... | 39 |
| II.3.1. Sơ đồ khối | 39 |
| II.3.2. Ảnh đầu vào | 39 |
| II.3.3. Module Tiền xử lý..... | 40 |
| II.3.4. Phân tích sử dụng giả pháp Fractal Signature | 41 |
| II.4. Kết luận | 45 |
| Chương III: XÂY DỰNG CHƯƠNG TRÌNH THỬ NGHIỆM..... | 46 |
| III.1. Yêu cầu hệ thống | 46 |
| III.2. Thiết kế chương trình..... | 46 |
| III.2.1. Cấu trúc dữ liệu..... | 46 |
| III.2.2. Module chuẩn hóa ảnh | 48 |
| III.2.3. Module giao diện chính..... | 51 |
| III.2.4. Module phân tích Top-down (TD)..... | 52 |
| III.2.5. Module phân tích Fractal Signature | 55 |
| III.2.6. Module lọc và làm trơn nhiễu..... | 57 |
| III.2.7. Module mô phỏng thuật toán FS | 58 |
| III.2.8. Các hàm chức năng chính của image processing tool trong matlab sử dụng trong chương trình | 60 |
| III.3. Kết luận và đánh giá kết quả | 62 |
| Kết luận..... | 83 |
| TÀI LIỆU THAM KHẢO | 84 |
| Phục Lục | 85 |
| A. Mã nguồn đầy đủ của chương trình | 85 |
| A.1. Danh mục các chương trình con trong chương trình | 85 |
| A.2. Sơ khối liên kết giữa các thủ tục trong chương trình..... | 86 |

A.3. Mã nguồn các module 86

Danh mục các hình ảnh

| | |
|---|----|
| Hình 1: Sơ đồ tổng quan quá trình tạo ảnh tài liệu | 5 |
| Hình 2: Ví dụ ảnh tài liệu | 6 |
| Hình 3: Sơ đồ OCR cơ bản | 7 |
| Hình 4: b-Cấu trúc vật lý; c,d-Cấu trúc logic của một tài liệu[4] | 9 |
| Hình 5: Ví dụ loại tài liệu có bố cục phức tạp | 10 |
| Hình 6: Sơ đồ nguyên lý hệ thống xử lý tài liệu[6] | 11 |
| Hình 7: a - Ảnh gốc b - Ảnh sau khi tách nền | 12 |
| Hình 8: Ví dụ một ảnh tài liệu bị nghiêng một góc 5 độ | 13 |
| Hình 9: Ví dụ một cây mô tả cấu trúc logic của một trang tài liệu[5] | 14 |
| Hình 10: VnDOCR và một ví dụ nhận dạng | 15 |
| Hình 11: Ảnh mẫu có cấu trúc vật lý phức tạp | 16 |
| Hình 12: Kết quả ra hai vùng ảnh với ảnh mẫu 11 | 16 |
| Hình 13: Mẫu ảnh có cấu trúc vật lý phức tạp, nhưng các khối bao bởi hình chữ nhật | 17 |
| Hình 14: Kết quả phân tích với ảnh 13 | 18 |
| Hình 15: Đầu ra phân vùng chỉ có 1 vùng văn bản | 19 |
| Hình 16: Đầu ra có vùng chứa cả ảnh và text | 19 |
| Hình 17: Với ảnh 13 đạt hiệu quả 90% | 20 |
| Hình 18 Với ảnh I-15 hiệu quả đạt 100% | 21 |
| Hình 19: Với mẫu phức tạp hơn Finereader cho kết quả 95% | 22 |
| Hình 20: Kết quả chiếu nghiêng theo phương ngang và phương thẳng đứng của một trang tài liệu | 24 |
| Hình 21: Phân tách cột dựa vào phép chiếu nghiêng theo phương ngang | 25 |
| Hình 22: Phép chiếu nghiêng theo phương ngang để phân đoạn ký tự hoặc từ | 26 |
| Hình 23: Lược đồ chiếu ngang của một dòng chữ nghiêng - rất khó phân đoạn ký tự | 27 |
| Hình 24: Lược đồ chiếu đứng của trang tài liệu bị nghiêng | 28 |
| Hình 25: Lược đồ chiếu đứng của một bài báo | 29 |
| Hình 26: Phương pháp Dostrum cho phân tích định dạng trang từ dưới lên. (a) Một phần của nội dung văn bản gốc. (b) Các thành phần lân cận gần nhất được xác định. | |

| | |
|---|----|
| (c) Các hình chữ nhật tối thiểu tạo nên nhóm láng giềng gần nhất từ đó xác định được dòng văn bản. | 31 |
| Hình 27: Mô tả thuật toán Tách và Nối thích nghi | 33 |
| Hình 28: Mô tả thuật toán FS..... | 35 |
| Hình 29: Sơ đồ khối hệ thống phân tích tài liệu trong phạm vi đề tài | 39 |
| Hình 30: Ví dụ một block chuyển sang dạng bề mặt trong không gian 3D | 41 |
| Hình 31: Ví dụ chuyển ảnh chữ "c" sang dạng bề mặt trong không gian 3D | 41 |
| Hình 32: (a) Ảnh một tài liệu gốc, (b) kết quả sau khi áp dụng FS | 44 |
| Hình 33: Giao diện chính..... | 51 |
| Hình 34: Kết quả phân tích của top-down trên tài liệu có cấu trúc đơn giản..... | 63 |
| Hình 35: Kết quả phân tích của top-down trên tài liệu có cấu trúc bảng..... | 65 |
| Hình 36: : Kết quả phân tích của top-down trên tài liệu thuần văn bản..... | 67 |
| Hình 37: : Kết quả phân tích của top-down trên tài liệu có cấu trúc phức tạp (trang tạp trí)..... | 69 |
| Hình 38 : Kết quả phân tích của top-down trên một tờ quảng cáo | 71 |
| Hình 39: Kết quả phân tích của top-down trên tài liệu có cấu trúc phức tạp..... | 73 |
| Hình 40: Kết quả phân tích của top-down trên tài liệu bị nghiêng..... | 75 |
| Hình 41: Kết quả phân tích của FS trên tài liệu bị nghiêng | 77 |
| Hình 42: Kết quả phân tích của FS trên tài có cấu trúc phức tạp | 79 |
| Hình 43: Kết quả phân tích của FS trên một trang quảng cáo..... | 81 |
| Hình 44: Kết quả phân tích của FS trên tài liệu đơn giản | 82 |

MỞ ĐẦU

I. Đặt vấn đề

Ngày nay việc sử dụng máy tính để lưu trữ tài liệu không còn là vấn đề mới mẻ và cần phải chứng minh tính an toàn, thuận tiện của nó. Tuy nhiên việc sử dụng giấy để lưu trữ tài liệu trong một số mục đích vẫn không thể thay thế được (như báo, sách, công văn,...). Hơn nữa lượng tài liệu được tạo ra từ nhiều năm trước vẫn còn rất nhiều mà không thể bỏ đi được vì tính quan trọng của chúng.

Chúng ta mong muốn có thể điện tử hóa hàng tỉ trang tài liệu đó và cất chúng chỉ trong một ổ cứng kích thước bằng một cuốn sách nhỏ, tìm kiếm thông tin mà chỉ cần tốn vài giây với một cái gõ phím Enter. Giải pháp là gì?

Thông thường người ta sẽ phải thuê người cùng với việc tốn hàng tháng, hàng năm mới có thể nhập vào máy tính được hết lượng tài liệu đó. Hiện nay chúng ta đã có các máy Scan với tốc độ cao, công nghệ xử lý của máy tính ngày càng siêu việt với tốc độ tính toán vượt cả tốc độ ánh sáng, vậy tại sao chúng ta không quét toàn bộ các trang tài liệu vào và chuyển chúng thành văn bản một cách tự động?

Bằng cách đó tốc độ và tính chính xác sẽ tăng hàng trăm lần trong khi chi phí lại là cực tiểu. Vấn đề là khi quét vào máy tính chúng ta không thu được ngay các dòng văn bản từ các trang tài liệu kia, để có thể soạn thảo, sửa chữa và tìm kiếm như làm trên Office. Tất cả những gì thu được chỉ là các tấm ảnh của các trang văn bản, máy tính lại đối xử công bằng như nhau với mọi điểm ảnh, máy tính không có “mắt” như chúng ta để biết đâu là điểm ảnh của chữ, đâu là điểm ảnh của đối tượng đồ họa.

Một giải pháp được nghĩ đến ngay đó là đó là xây dựng các hệ thống nhận dạng chữ, trong tấm ảnh chứa cả chữ và đối tượng đồ họa cần tách và chuyển thành dạng trang văn bản, từ đó có thể mở và soạn thảo được trên các trình soạn thảo văn bản.

Một cách tổng quát thì cách thức làm việc của một hệ thống nhận dạng chữ như sau[5]:

1. Chụp ảnh các trang tài liệu trên giấy và lưu lại trong máy tính dưới dạng hình ảnh.

2. Sử dụng một chương trình xử lý ảnh để phân tích hình ảnh sau khi quét, đọc được ký tự trên hình ảnh đó và ghi lại vào máy tính theo cách mà máy tính quản lý được thông tin dữ liệu đó.
 - a. Bước 1 là phân tích cấu trúc của ảnh tài liệu, từ đó xác định đâu là phần chứa chữ, đâu là phần chứa cả ảnh lẫn ký tự và đâu chỉ chứa hình ảnh. Bước này thực sự quan trọng cho bước nhận dạng. Bởi nó định vị chính xác cho việc áp dụng các thuật toán nhận dạng lên vùng đã xác định tính chất, nếu bước này chính xác trước tiên nó hạn chế thời gian cho việc nhận dạng, sau là tăng ngữ nghĩa bổ sung cho việc nhận dạng.
 - b. Bước 2 nhận dạng ký tự dựa vào các tính chất của ký tự, ví dụ như sắp xếp theo dòng, khoảng cách giữa 2 từ lớn hơn khoảng cách giữa 2 ký tự, dùng trí tuệ nhân tạo để dự đoán các ký tự kề nhau phải như thế nào, các từ trong câu phải như thế nào để câu có nghĩa. Từ đó có nội dung đúng để lưu trữ, quản lý....

Trong thực tế không phải quá trình nhận dạng nào cũng chỉ trải qua hai bước như trên, bởi vì có rất nhiều tham số ảnh hưởng đến kết quả của các chương trình nhận dạng, như nhiều, Font chữ, kích thước chữ, kiểu chữ nghiêng, đậm, gạch dưới. Ngoài ra các dòng chữ cũng có thể trộn lẫn với các đối tượng đồ họa, vì thế trước khi nhận dạng chữ, một số thao tác tiền xử lý sẽ được tác động lên ảnh như, lọc nhiễu, chỉnh góc nghiêng và đặc biệt quan trọng là phân tích trang tài liệu để xác định cấu trúc của trang văn bản đồng thời tách biệt hai thành phần là chữ và các đối tượng đồ họa (phi chữ).

II. Nội dung nghiên cứu

1. Mục tiêu nghiên cứu chính của đề tài

- Tìm hiểu cấu trúc trang tài liệu (cấu trúc vật lý, logic)?
- Tìm hiểu một số kỹ thuật phân tích trang tài liệu (phân vùng, phân đoạn,...)
- Cài đặt thử nghiệm một giải pháp phân tích có hiệu quả cao so với các phương pháp truyền thống như top-down hay bottom-up trên **ảnh vào là ảnh đa cấp xám có cấu trúc phức tạp**.
- Từ kết quả nghiên cứu có một sự chuẩn bị kiến thức đầy đủ cho bước nghiên cứu tiếp theo là nhận dạng ký tự quang.

2. Ý nghĩa khoa học của đề tài

- Giải quyết được vấn đề về học thuật: đề tài sẽ mang ý nghĩa cung cấp về mặt lý thuyết để làm rõ về các phương pháp phân tích trang tài liệu.
- Đáp ứng được yêu cầu của thực tiễn: từ các lý thuyết đã được nghiên cứu, từ đó liên hệ và gắn vào thực tiễn để có thể áp dụng vào các lĩnh vực như: Lưu trữ thư viện, điện tử hóa văn phòng, nhận dạng và xử lý ảnh, ...

3. Nhiệm vụ nghiên cứu

Mục đích của luận văn đề cập được đến hai phần:

- Phần lý thuyết: Nắm rõ và trình bày những cơ sở lý thuyết liên quan đến cấu trúc trang tài liệu, một số kỹ thuật phân tích trang tài liệu, từ đó có thể xác định tính quan trọng của bước này trong nhận dạng ký tự, đồng thời hiểu các công việc kế tiếp cần làm trong bước nhận dạng ký tự.
- Phần phát triển ứng dụng: Áp dụng các thuật toán đã trình bày ở phần lý thuyết từ đó lựa chọn một giải pháp tối ưu và cài đặt thử nghiệm chương trình phân tích trang tài liệu.

4. Phương pháp nghiên cứu

- Tìm kiếm, tham khảo, tổng hợp tài liệu từ các nguồn khác nhau để xây dựng phần lý thuyết cho luận văn.
- Sử dụng các kỹ thuật được áp dụng phân tích trang tài liệu để làm rõ bản chất của các vấn đề được đưa ra trong phần lý thuyết.
- Xây dựng chương trình Demo.

5. Phạm vi nghiên cứu

Bài toán nhận dạng và xử lý ảnh tài liệu đã được phát triển với nhiều thành tựu trong thực tế, có rất nhiều thuật toán tối ưu đã được các nhà khoa học đề nghị. Tuy nhiên có thể nói chưa có một chương trình nào có thể “đọc” một ảnh văn bản như con người, vì thực tế có rất nhiều kiểu trang văn bản khác nhau, khác nhau về cấu trúc trình bày, ngôn ngữ, kiểu font, chữ viết tay,... Đây thực sự là một bài toán lớn, chính vì thế trong phạm vi của luận văn chỉ tìm hiểu một số kỹ thuật phân tích trang văn bản tiêu biểu với mục đích để so sánh và một thuật toán mới chưa được đưa ra ở các đề tài trước. Cuối cùng, dựa vào đó để xây dựng Demo cho một ứng dụng.

Các kết quả nghiên cứu dự kiến cần đạt được:

- Tìm hiểu tài liệu liên quan đến lĩnh vực quan tâm để nắm bắt được bản chất vấn đề đặt ra.
- Báo cáo lý thuyết
- Chương trình Demo.

III. Bộ cục của luận văn

Nội dung của luận văn được trình bày trong ba chương với nội dung chính sau.

Chương 1: Trình bày các khái niệm và mô hình tổng quát của hệ thống nhận dạng chữ viết, cùng với một số phần mềm nhận dạng tiêu biểu hiện nay.

Chương 2: Trình bày một số phương pháp phân tích trang tài liệu, từ đó đánh giá ưu nhược điểm để lựa chọn phương pháp Fractal Signature cho chương trình thử nghiệm. Trình bày về thiết kế cho chương trình demo.

Chương 3: Trình bày chi tiết về việc cài đặt chương trình cũng như các thủ tục sử dụng trong chương trình với phương pháp phân tích Fractal Signature và ảnh đầu vào là ảnh đa cấp xám có độ phức tạp cao.

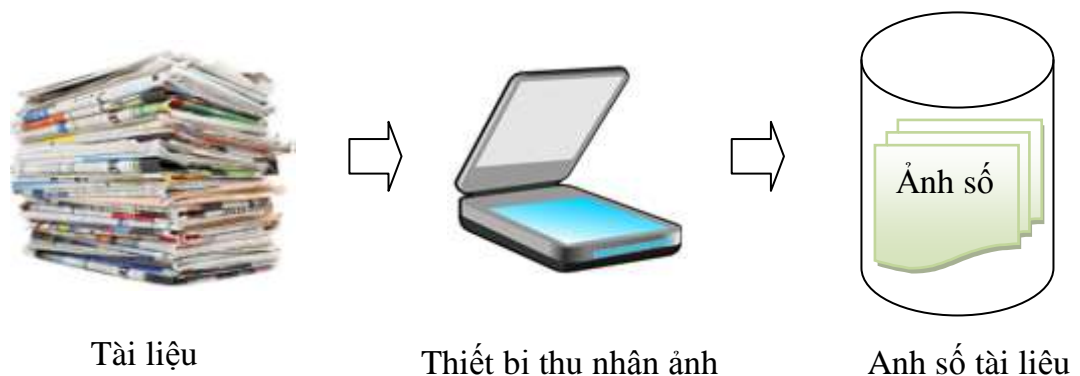
Chương I. TỔNG QUAN VỀ NHẬN DẠNG CHỮ VIẾT VÀ PHÂN TÍCH TRANG TÀI LIỆU

Chương này đưa ra các khái niệm về đối tượng làm việc của đề tài là ảnh tài liệu, khái niệm về cấu trúc vật lý và cấu trúc logic. Giới thiệu các khâu trong một hệ thống nhận dạng chữ viết hoàn chỉnh. Đồng thời đưa ra một số phần mềm nhận dạng của Việt Nam và Thế giới cùng với các mẫu kết quả phân tích của nó nhằm mục đích so sánh và xác định phạm vi cho đề tài.

I.1. Ảnh tài liệu và nhận dạng ảnh tài liệu

I.1.1. Tổng quan về ảnh tài liệu

Trang ảnh tài liệu được đề cập ở đây là các file ảnh số hoá thu được bằng cách quét các trang tài liệu dùng máy scanner, máy ảnh số, hay nhận từ một máy fax (Hình 1), file ảnh này được lưu giữ trong máy tính. Ảnh tài liệu có nhiều loại: ảnh đen trắng, ảnh màu, ảnh đa cấp xám với các phần mở rộng như TIF, BMP, PCX, ... (Hình 2) và ảnh tài liệu được đưa ra trong luận văn này là **ảnh đa cấp xám**.



Hình 1: Sơ đồ tổng quan quá trình tạo ảnh tài liệu