

ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH
TRƯỜNG ĐẠI HỌC BÁCH KHOA
KHOA KHOA HỌC - KỸ THUẬT MÁY TÍNH



LUẬN VĂN TỐT NGHIỆP ĐẠI HỌC

Nhận dạng biểu thức toán học viết tay

Giảng viên hướng dẫn : TS. Lê Thành Sách

Nhóm sinh viên thực hiện : Phan Tấn Phúc - 51303058

Bùi Khánh Ngọc - 51302567

Tp. Hồ Chí Minh, Tháng 05/2017



Lời nói đầu

Thực tập tốt nghiệp là giai đoạn báo hiệu một chặng đường Đại học đã sắp kết thúc. Đây cũng là giai đoạn chuẩn bị quan trọng cho sự thành công của Luận văn sau này. Để đi được tới thời điểm này, nhóm muốn gửi lời cảm ơn chân thành và sự biết ơn đến gia đình- nguồn động lực cũng như nguồn hỗ trợ tài chính cho các thành viên nhóm có điều kiện tốt để học tập, cảm ơn các bạn đã cùng lên lớp, cùng làm bài tập, cùng chơi vui cũng như đã giúp đỡ nhóm trong các học kỳ vừa qua và không thể không nhắc đến các thầy cô trong Khoa- những người đã cho chúng em tiếp cận những kiến thức mới và dạy chúng em cách trưởng thành.

Ngoài ra, nhóm muốn gửi lời cảm ơn đặc biệt đến thầy hướng dẫn- Tiến sĩ Lê Thành Sách. Cảm ơn thầy đã cho nhóm cơ hội được làm việc với thầy, cảm ơn thầy vì đã luôn theo sát, hỗ trợ cũng như định hướng trong công việc cho nhóm và cảm ơn thầy về những bài học trên phương diện làm người.

Sau cùng, vì những hạn chế về mặt thời gian cũng như khả năng trong cách trình bày và viết báo cáo nên không thể tránh khỏi những thiếu sót, rất mong sự thông cảm và những ý kiến góp ý từ quý thầy cô và các bạn để giúp nhóm hoàn thiện hơn trong giai đoạn sau này.

Chân thành cảm ơn.



Tóm tắt báo cáo

Báo cáo này chỉ ra quá trình thực hiện của nhóm để bước đầu tiếp cận với đề tài cũng như những kiến thức có được từ quá trình đó. Từ đó tạo cơ sở cho giai đoạn luận văn sau này. Bố cục của bài báo cáo bao gồm 4 chương, không kể các mục lục, phụ lục khác.

Chương 1 là chương giới thiệu tổng quan đề tài. Ở đó sẽ trình bày thế nào là nhận dạng biểu thức toán học, nhu cầu của nó trong xã hội cũng như đưa ra lý do vì sao nhóm chọn đề tài này và quá trình thực hiện của nhóm cho đến thời điểm hiện tại.

Chương 2 sẽ trình bày các kiến thức mà nhóm tìm hiểu được liên quan đến đề tài. Đó là những kiến thức cơ bản về mạng nơron nhân tạo dựa trên phân tích mạng truyền thẳng 1 lớp, mạng nơron tích chập cùng với một kiến trúc nổi tiếng của nó là Lenet-5.

Chương 3 là tóm lược những điều đã đọc được từ các bài báo viết về những công trình liên quan trực tiếp đến đề tài.

Chương 4 là chương cuối cùng với nhiệm vụ là tổng kết, đánh giá những việc nhóm đã làm được cho đến hiện tại và kế hoạch phát triển cho giai đoạn tới.

Mục lục

Mục lục	3
Danh sách hình vẽ	5
Chương 1 Giới thiệu	7
1 Giới thiệu đề tài	7
2 Lý do chọn đề tài	7
3 Phạm vi đề tài	8
4 Quá trình thực hiện	8
Chương 2 Kiến thức đã tìm hiểu	10
1 Mạng nơ-ron tích chập (CNN)	10
2 Mạng Single Shot Multibox Detector (SSD)	10
2.1 Bộ mã hóa (Encoder)	10
2.2 Bộ phát hiện - phân loại	12
2.3 Tính giá trị mất mát	13
2.4 Bộ giải mã (Decoder)	15
2.5 Một số vấn đề khác trong mạng SSD	15
Chương 3 Công trình liên quan	16
1 Cái nhìn toàn cảnh	16
2 WHOOOOO : Watch, Attend and Parse: An End-to-end Neural Network Based Approach to Handwritten Mathematical Expression Recognition . .	17
2.1 Watcher	17
2.2 Attend	17
2.3 Parser	17
3 Wenhao He[1]: Context-aware Recognition	17
4 QAK[2]	19
Chương 4 Mô hình đề xuất	20
1 Tổng quan	20
2 Phát hiện ký tự	20
2.1 Mạng cơ sở	20
2.2 Thân mạng SSD	20
2.3 Các thay đổi	23
3 Phân tích cú pháp	25
3.1 Sinh cây BST	25
3.2 Sinh cây Lexed - BST	29
Chương 5 Hiện thực, đánh giá	30
1 Chuẩn bị dữ liệu	30



2	Hiện thực hệ thống	30
2.1	Quá trình tiền huấn luyện	30
3	Xây dựng bản thử nghiệm	31
4	Đánh giá ưu, nhược điểm	32
4.1	Ưu điểm	32
4.1.a	Nhược điểm	32
Chương 6	Tổng kết	32
1	Công việc đã làm được	32
1.1	Chuẩn bị cơ sở lý thuyết	32
2	Nhận xét bản thử nghiệm hiện tại	32
3	Kế hoạch trong giai đoạn tới	33
Chương 7	Kết luận	33
1	Kết luận	33
2	Hướng phát triển trong tương lai	33
Tài liệu		34

Danh sách hình vẽ

1	Cấu tạo của một mạng SSD	13
2	Mô hình học được đề xuất [1].	18
3	Quy trình nhận dạng.	19
4	Cấu tạo mạng VGG16	20
5	Sơ đồ các lớp trong mạng SSD	21
6	Ảnh ví dụ chữ nhỏ (với kích thước thật)	23
7	Giao diện bản thử nghiệm	31



Danh mục từ viết tắt

Thuật ngữ	Giải thích
CNN	Mạng nơron tích chập (Convolutional Neural Network)
NN	Mạng nơron truyền thống
MSE	Mean Square Error



Chương 1 Giới thiệu

1 Giới thiệu đề tài

Sự phát triển của khoa học công nghệ cùng với sự bùng nổ của thiết bị di động thúc đẩy một cuộc "cách mạng trên thiết bị cầm tay". Con người đòi hỏi nhiều hơn ngoài các tính năng nghe, gọi, chụp ảnh, giải trí thông thường trên điện thoại. Thực tế đã có rất nhiều ứng dụng điện thoại nói riêng và thiết bị di động nói chung ra đời đáp ứng nhu cầu ngày càng cao của xã hội. Trong lĩnh vực y tế, phải kể đến ứng dụng giám sát sức khỏe người dùng trên những chiếc đồng hồ thông minh của Apple hay Samsung. Với giao thông, những ứng dụng chỉ đường, định vị và giám sát xe đề phòng trộm ngày càng phổ biến và giúp ích thực sự cho con người. Riêng với giáo dục, vấn đề thường gặp đối với các bạn học sinh là giải, trực quan hoá các phương trình, hàm số toán học hay soạn các giáo án chứa nhiều công thức, ký hiệu phức tạp đối với thầy cô. Họ cần một giải pháp nào đó giúp giảm thiểu công sức trong những tình huống như vậy. Giải pháp này có thể giải quyết những vấn đề cơ bản sau:

- Số hoá các công thức in trong sách hay được viết tay.
- Giải tham khảo một số dạng phương trình.
- Biểu diễn công thức, ký hiệu dưới những lệnh mà các trình soạn thảo toán có thể hiểu được, ví dụ LaTeX.
- ...

2 Lý do chọn đề tài

Bởi sự cần thiết về một ứng dụng nhận dạng biểu thức toán học đã được trình bày ở mục 1, trên thị trường cũng đã xuất hiện nhiều sản phẩm như vậy, đáng chú ý là PhotoMath¹. Tuy nhiên, nhóm nhận thấy thách thức nếu phải hiện thực thành công đề tài này. Một số câu hỏi đã được đặt ra:

- Làm sao có thể nhận dạng được các ký hiệu?
- Làm cách nào để nhận dạng cả một biểu thức?
- Làm sao biết được đây là loại biểu thức gì?
- Cách viết thì khác nhau với từng người sẽ ảnh hưởng đến kết quả nhận dạng, vậy có cách nào để khắc phục?

¹Một ứng dụng về nhận dạng và giải các biểu thức toán học nổi bật trên Google Play.



- Nhóm có thể tạo ra được một sản phẩm hoàn thiện như PhotoMath không?

Để tự mình trả lời những câu hỏi đó, nhóm quyết tâm thực hiện đề tài này. Ngoài ra, việc áp dụng kiến thức đã học để tạo ra một sản phẩm vừa cần thiết cho xã hội vừa tự bản thân mình có thể sử dụng được tạo cho nhóm một động lực để tiến hành.

3 Phạm vi đề tài

- Nhận dạng biểu thức toán học viết tay dạng offline².
- Chuyển biểu thức từ dạng hình ảnh sang dạng máy có thể hiểu được (Latex).

4 Quá trình thực hiện

Bước 1: Tìm hiểu công trình đã được hiện thực bởi nhóm sinh viên khoá 2011 [2].³

- Mục tiêu:
 - Hình dung ban đầu về đề tài.
 - Biết được các kiến thức mà nhóm sinh viên khoá 2011 đã sử dụng để giải quyết đề tài.
- Công việc:
 - Liên hệ nhóm sinh viên khoá 2011 để có được toàn bộ mã nguồn cũng như tài liệu liên quan.
 - Đọc tài liệu và chạy lại bộ mã nguồn.

Bước 2: Tìm hiểu các kiến thức nền liên quan đến xử lý ảnh và nhận dạng.

- Mục tiêu:
 - Hiểu được các kiến thức cơ sở phục vụ đề tài.
 - Tạo điều kiện thuận lợi cho việc đọc các bài báo liên quan đề tài.
- Công việc:
 - Tìm đọc và hiểu rõ kiến thức được giới thiệu bởi thầy hướng dẫn.
 - Mở rộng kiến thức bằng cách tìm kiếm các kiến thức liên quan.
 - Hiện thực một số kiến thức để hiểu rõ hơn.

²Nhận dạng từ ảnh chứa biểu thức toán học

³Đề tài nhận dạng biểu thức toán học đã từng được hiện thực bởi một nhóm sinh viên khoá 2011.



Bước 3: Tìm hiểu công trình nghiên cứu gần đây nhất liên quan đến đề tài.

- Mục tiêu:
 - Hiểu rõ bản chất của đề tài cũng như những hướng phát triển có thể có.
 - Biết được phương pháp mà các nhóm tác giả sử dụng cho công trình của mình.
 - Có được những đánh giá, nhận xét ban đầu về các phương pháp đã được sử dụng.
- Công việc:
 - Tìm và đọc các bài báo khoa học gần nhất đề cập đến đề tài.
 - Tìm hiểu những kiến thức được nêu ra trong bài báo.

Bước 4: Đề xuất phương pháp.

- Mục tiêu:
 - Có được chương trình thử nghiệm nhận dạng biểu thức toán học.
- Công việc:
 - Hiện thực đề tài theo phương pháp được nêu ra bởi nhóm sinh viên khoá 2011.

Bước 5: Đánh giá kết quả thực hiện

- Mục tiêu:
 - Có được đánh giá ban đầu và điều chỉnh phương pháp cho phù hợp.
- Công việc:
 - Xem xét ưu, nhược điểm của phương pháp.
 - Đề ra phương pháp cải tiến hoặc thay thế.

Chương 2 Kiến thức đã tìm hiểu

1 Mạng nơ-ron tích chập (CNN)

2 Mạng Single Shot Multibox Detector (SSD)

SSD như là một mạng cải tiến phương pháp phát hiện bằng cửa sổ trượt⁴. Thay vì sử dụng một (một số) cửa sổ có kích thước cố định, thì SSD sinh ra một số lượng hữu hạn các ô chuẩn⁵ để rời từ các ô chuẩn đó để hệ thống có thể xác định vị trí các ký tự cần nhận diện cho quá trình huấn luyện, qua đó mạng cần phải học cách dự đoán cả kích thước của các ô bọc quanh ký tự thay vì chỉ chấp nhận kích thước cho trước. SSD cũng sử dụng các lớp tích chập để trích đặc trưng, tạo tiền đề cho việc phát hiện và phân loại ký tự.

Nhóm xin phép được tách quá trình huấn luyện thành ba giai đoạn: mã hóa⁶, phát hiện - phân loại và tính giá trị mất mát. Và song song với huấn luyện, quá trình kiểm tra, kiểm định, chạy thực tiễn cũng được chia thành ba giai đoạn: trích đặc trưng - phân loại và giải mã⁷

2.1 Bộ mã hóa (Encoder)

Bộ mã hóa chỉ được sử dụng trong quá trình huấn luyện nhằm mã hóa, chuyển đổi từ "đáp án"⁸ thô (được lưu trong tệp tin txt hoặc xml) sang "đáp án" mà mạng SSD có thể hiểu được.

Sinh ô chuẩn

Để mã hóa "đáp án", bộ mã hóa trước hết có nhiệm vụ sinh ra nhiều ô chuẩn có nhiều kích thước khác nhau, để làm được điều này, mạng cần phải được cung cấp một số thông số:

- Danh sách kích thước của các feature map (Phần này sẽ được giải thích rõ hơn ở mục sau): Mỗi feature map trong danh sách tương ứng với một mức kích thước⁹ trong việc phát hiện ảnh. Trong một mức kích thước, tất cả các ô chuẩn đều có kích thước như nhau và cách đều nhau, mỗi pixel trong feature map thể hiện một (một số ô chuẩn), vì vậy, ta có thể tính được danh sách khoảng cách giữa trọng tâm giữa các ô chuẩn bằng cách lấy kích thước ảnh chia cho kích thước của feature map ở mức kích thước tương ứng. Tất cả các dữ liệu liên quan đến kích thước sau đó sẽ được chuyển về tập giá trị $\{x \in R | x \in [0, 1]\}$

⁴Thuật ngữ tiếng Anh: Sliding Window

⁵Thuật ngữ tiếng Anh: Default box

⁶Thuật ngữ tiếng Anh: Encoder

⁷Thuật ngữ tiếng Anh: Decoder

⁸Thuật ngữ tiếng Anh: Ground truth

⁹Thuật ngữ tiếng Anh: Scale

- Danh sách các tỉ lệ diện mạo¹⁰ ứng với mỗi mức kích thước. Mạng SSD dựa vào các tỉ lệ diện mạo đó để tạo ra các ô chuẩn có hình dạng khác nhau tại cùng một vị trí (trọng tâm của các ô chuẩn có cùng một vị trí).
- Kích thước nhỏ nhất và lớn nhất của các ô chuẩn.

Từ những thông số trên, các ô chuẩn sẽ được sinh ra theo các bước sau:
Trước hết, bộ phận này tiến hành sinh ra kích thước của ô chuẩn ứng với mỗi mức kích thước tương ứng bằng công thức:

$$s_k = s_{min} + \frac{s_{max} - s_{min}}{m - 1}(k - 1)$$

Trong đó s_{min} và s_{max} là kích thước nhỏ nhất và lớn nhất của ô chuẩn, m là số lượng feature map. Việc chọn các thông số này khá quan trọng vì nó sẽ ảnh hưởng đến các ký tự có thể nhận diện được sau này, các kích thước cần được chọn sao cho không quá lớn cũng như không quá nhỏ đối với những đối tượng được kỳ vọng nhận diện được.

Ứng với mỗi mức kích thước:

- Mỗi ô chuẩn đều được liên kết với một pixel trên feature map, ta cũng đã tính được danh sách khoảng cách trọng tâm theo đề cập ở trên. Từ đó ta dễ dàng tính được tọa độ của các trọng tâm của các ô chuẩn trong ảnh theo công thức:

$$(x, y) = \left(\left(i + \frac{1}{2} \right) \times step, \left(j + \frac{1}{2} \right) \times step \right)$$

Trong đó: x, y là tọa độ của các trọng tâm, i, j là các số nguyên dương nằm trong khoảng từ 0 đến kích thước feature map đang xét và $step$ là khoảng cách giữa hai trọng tâm liền kề trong mức kích thước đang xét.

- Sau khi có được tọa độ của các trọng tâm, bộ phận mã hóa tiến hành sinh các ô chuẩn ứng với mỗi vị trí trọng tâm:
 - Tạo một ô chuẩn có kích thước s_k là kích thước ứng với mức kích thước k hiện tại.
 - Tạo một ô chuẩn có kích thước bằng $\sqrt{s_k \times s_{k+1}}$.
 - Với mỗi phần tử trong danh sách tỉ lệ diện mạo tương ứng, ta tạo một ô chuẩn có kích thước chiều rộng: $w = s_k \times \sqrt{aspect_ratio}$ và chiều cao bằng $h = s_k \div \sqrt{aspect_ratio}$ với $aspect_ratio$ là tỉ lệ diện mạo đang xét.

Như vậy, với mỗi vị trí được chọn để đặt trong tâm các ô chuẩn, bộ phận mã hóa sẽ sinh ra $2 + n$ ô chuẩn với n là số tỉ lệ diện mạo ứng với mức kích thước hiện tại. Kết quả của quá trình sinh ô chuẩn là một tập hợp nhiều vector. Mỗi vector ứng với một ô chuẩn mang thông tin tọa độ trọng tâm, chiều dài, chiều rộng của ô chuẩn.

¹⁰Thuật ngữ tiếng Anh: Aspect ratio

Kết hợp¹¹

Sau khi hoàn tất sinh ô chuẩn, bộ phận encoder có nhiệm vụ kết hợp dữ liệu thô là các ô bọc¹² từ "đáp án"¹³ qua các ô chuẩn vừa được tạo. Kết quả là ta sẽ thu được một tập nhiều vector, mỗi vector đại diện cho mỗi ô chuẩn chứa thông tin gồm tọa độ trọng tâm, kích thước của ô chuẩn và nhãn của ô chuẩn đó là gì (có thể là phần nền - không có gì cả, hoặc một ký tự nào đó cần được nhận diện).

Việc kết hợp được thực hiện bằng cách so sánh từng ô chuẩn với từng ô bọc trong "đáp án" bằng cách tính chỉ số Jackard, nếu chỉ số này thỏa điều kiện thì ô chuẩn đó được gán nội dung là ký tự chứa trong ô bọc phù hợp nhất, nếu không thì ô chuẩn được gán nội dung là "nền".

Chỉ số Jackard giữa hai ô bọc được tính bằng công thức:

$$J(A, B) = \frac{A \cap B}{A \cup B}$$

Trong đó, A, B lần lượt là phần ảnh mà ô bọc A và ô bọc B chiếm. Công thức này mang ý nghĩa hai ô bọc có kích thước càng tương đồng và phần trùng nhau càng nhiều thì có chỉ số Jackard càng gần giá trị 1.

Các ô bọc chỉ là thành phần trung gian để mạng SSD có thể tính được giá trị mất mát cũng như thực hiện quá trình lan truyền ngược, vì vậy để thuận tiện cho việc tính mất mát sau này, từ dữ liệu vị trí và kích thước của ô bọc, bộ mã hóa tạo ra các vector có các giá trị biểu diễn sự mất mát của các ô bọc so với các ô chuẩn, sau đó nhiệm vụ của mạng là dự đoán ra sự mất mát đó (Vấn đề này sẽ được đề cập rõ hơn trong [hần tính giá trị mất mát]).

2.2 Bộ phát hiện - phân loại

Đây là bộ phận có nhiệm vụ trích đặc trưng từ ảnh thô đầu vào, tạo ra các feature map để đưa vào quá trình phân loại ký tự.

Trích đặc trưng (Hay mạng cơ sở¹⁴)

Mạng cơ sở của SSD là một mạng nơron tích chập bất kỳ, có thể là VGG[3] hoặc lenet[4], ... Hầu hết những mạng này có mục đích phân loại¹⁵, vì vậy ở phía cuối mạng thường có các lớp liên kết đầy đủ nhằm giảm kích thước feature map và tạo ra dữ liệu đầu ra có kích thước phù hợp (bằng với số nhân cần dự đoán). Khi kết hợp với mạng SSD, các lớp liên

¹¹Thuật ngữ tiếng Anh: Matching

¹²Thuật ngữ tiếng Anh: Bounding box

¹³Thuật ngữ tiếng Anh: Ground truth

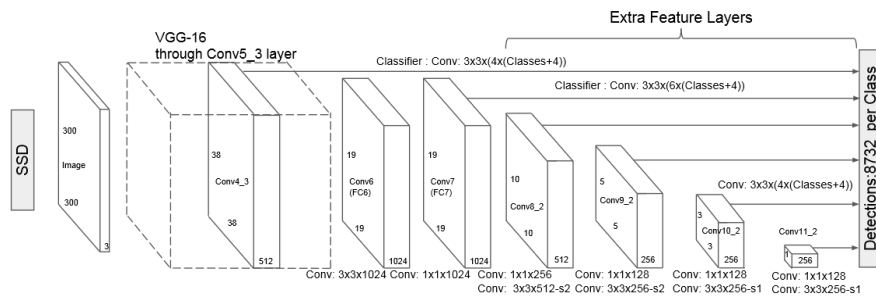
¹⁴Thuật ngữ tiếng Anh: Base Network

¹⁵Thuật ngữ tiếng Anh: Classification

kết đầy đủ này sẽ được thay thế bằng các lớp tích chập và song hành là các bước trích ra feature map để đưa vào lớp Multibox, chi tiết phần này sẽ được đề cập ở mục tiếp theo.

Phân loại

Bộ phận phân loại có nhiệm vụ đưa ra dự đoán về vị trí và nhãn của các ký tự có trong ảnh. Để làm được điều đó, bộ phân loại phải sinh ra các feature map ứng với từng mức kích thước khác nhau, các feature map đó được đưa vào lớp Multibox để sinh ra các vị trí dự đoán và nhãn tương ứng. Các feature map được lấy sau một (một số) lớp tích chập. Danh sách feature map có được ở giai đoạn sinh ô chuẩn chính là được lấy từ đây.



Hình 1: Cấu tạo của một mạng SSD

Sau khi có được các feature map, lớp Multibox sẽ tiến hành đưa ra dự đoán về ô bọc và nhãn gắn với ô bọc đó. Mỗi feature map sử dụng một tập hợp các lớp tích chập sẽ cho ra một số lượng dự đoán nhất định. Với mỗi vị trí trên feature map được kernel trượt qua, thì một dự đoán về ký tự được sinh ra. Kernel có kích thước:

$$3 \times 3 \times (n \times (classes + 4))$$

Với n là số lượng tỉ lệ diện mạo với fearture map tương ứng và $classes$ là số lượng nhãn cần nhận diện, số 4 trong công thức đại diện cho 4 thông số về vị trí của ô bọc (tọa độ trọng tâm và kích thước). Do kích thước của các kernel không đồng nhất (vì số lượng tỉ lệ diện mạo ứng với từng feature map có thể khác nhau), nên để tạo ra dữ liệu đầu ra phù hợp với các ô chuẩn đã tạo từ trước thì các vector dự đoán sẽ được sắp xếp lại sao cho dữ liệu đầu ra chỉ có $classes + 4$ kênh và vị trí của các vector phải tương ứng với vị trí các ô chuẩn đã sinh từ trước.

2.3 Tính giá trị mất mát

Giá trị mất mát được tính bằng tổng có trọng số giữa mất mát về vị trí¹⁶ và mất mát về độ tin cậy¹⁷

¹⁶Thuật ngữ tiếng Anh: Localization Loss

¹⁷Thuật ngữ tiếng Anh: Confidence Loss

$$L(x, c, l, g) = \frac{1}{N}(L_{conf}(x, c) + \alpha L_{loc}(x, l, g))$$

Trong đó:

- x biểu thị các phép kết hợp, cụ thể x_{ij}^p biểu thị phép kết hợp giữa ô chuẩn thứ i với ô bọc thứ j đối với nhân p . Vì vậy x_{ij}^p có tập giá trị $\{0, 1\}$.
- c biểu thị nhãn kỳ vọng cho phép kết hợp đang xét.
- l biểu thị ô bọc dự đoán.
- g biểu thị ô bọc "đáp án".

Hàm mất mát về vị trí có thể được biểu diễn bằng công thức:

$$L_{loc}(x, l, g) = \sum_{i \in Pos}^N \sum_{m \in \{cx, cy, w, g\}} x_{ij}^k smooth_{L1}(l_i^m - \hat{g}_j^m)$$

Trong đó:

- N là số lượng phép kết hợp giữa ô chuẩn và ô bọc có ý nghĩa (nhãn của phép kết hợp không phải là "nền". Nếu như $N = 0$ thì ta đặt giá trị mất mát bằng 0.
- Pos là tập hợp các phép kết hợp có ý nghĩa.
- $smooth_{L1}$ là một hàm tính mất mát được định nghĩa là

$$L(x) = \begin{cases} 0.5x^2 & \text{if } |x| < 1 \\ |x| - 0.5 & \text{otherwise} \end{cases}$$

- $\hat{g}_j^{cx} = (g_j^{cx} - d_i^{cx})/d_i^w$ Với d biểu thị cho ô chuẩn.
- $\hat{g}_j^{cy} = (g_j^{cy} - d_i^{cy})/d_i^h$
- $\hat{g}_j^w = \log\left(\frac{g_j^w}{d_i^w}\right)$
- $\hat{g}_j^h = \log\left(\frac{g_j^h}{d_i^h}\right)$

Vậy ta có thể thấy khi tính giá trị mất mát về vị trí, ta tính dựa trên sai lệch so với ô chuẩn mà ô bọc "đáp án" ban đầu kết hợp được thay vì tính mất mát trực tiếp.

hàm mất mát về độ tin cậy có thể được tính theo công thức:

$$L_{conf}(x, c) = - \sum_{i \in Pos}^N x_{ij}^p \log(\hat{c}_i^p) - \sum_{i \in Neg} \log(\hat{c}_i^0)$$

$$\text{Với } \hat{c}_i^p = \frac{\exp c_i^p}{\sum_p \exp c_i^p}$$



2.4 Bộ giải mã (Decoder)

Bản chất của bộ giải mã khá đơn giản, nó có chức năng chuyển dữ liệu các ô bọc mà mạng dự đoán (tức là những ô bọc được tính kích thước và tọa độ theo ô chuẩn mà ô bọc đó được kết hợp vào) sang dữ liệu tọa độ của ảnh (có gốc tọa độ nằm ở góc trên bên phải và tọa độ nằm trong miền $[0, 1]$). Dữ liệu này cần thiết để trực quan hóa các ký tự mà hệ thống dự đoán và để đưa vào bộ phân tích cú pháp để sinh mã Latex.

2.5 Một số vấn đề khác trong mạng SSD

Tăng cường dữ liệu¹⁸

Ngay trước quá trình kết hợp, ảnh và các ô bọc được đi qua một bước gọi là cắt ngẫu nhiên¹⁹. Ở bước này, ảnh được cắt đi một số phần ngẫu nhiên (có tác dụng giống như phóng to (zoom) ảnh. Điều này giúp làm đa dạng dữ liệu huấn luyện và làm giảm bớt hiện tượng overfit.

Cân bằng nhãn²⁰

Ngay sau bước kết hợp, hầu hết các ô

¹⁸Data Augumentation

¹⁹Thuật ngữ tiếng Anh: Random Crop

²⁰Hard Negative Mining

Chương 3 Công trình liên quan

Là một phần quan trọng của hệ thống **Nhận dạng ký tự thuộc thị giác**²¹, nhận dạng **biểu thức toán học**²² đã được nghiên cứu trong hơn nửa thế kỷ qua. Trên hành trình đó, rất nhiều công trình nhằm giải quyết vấn đề này đã được công bố. Trong chương này, nhóm sẽ trình bày tóm lược về một số công trình tiêu biểu, là những điểm tham khảo cho hệ thống nhận dạng biểu thức toán học mà nhóm sẽ hiện thực sau này và nhằm giúp cho bạn có một cái nhìn toàn cảnh về lĩnh vực này.

1 Cái nhìn toàn cảnh

Trong suốt hơn 50 năm hành trình, rất nhiều phương pháp nhận diện biểu thức toán học khác nhau đã được đề xuất, tuy vậy, bài toán này có thể được phân ra thành một số hướng tiếp cận khác nhau:

- Nhận diện theo hướng phân vùng²³. Ở hướng này, bài toán được chia ra thành nhiều bài toán nhỏ hơn thông qua việc phân chia vùng ảnh biểu thức toán học một cách có quy tắc. Một số ví dụ cho phương pháp này như giải thuật X-Y cut [5] hay dùng phương pháp projection profiles [6]. Đối với phương pháp này, bài toán sẽ trở nên vô cùng khó khăn với những ký tự dính liền nhau hoặc với những bài toán có dấu căn.
- Nhận diện dựa trên cấu trúc cây hoặc đồ thị. Ở hướng này, dựa vào vị trí các ký tự cũng như cấu trúc tổng thể của biểu thức, biểu thức toán học được biểu diễn theo cấu trúc cây hoặc đồ thị. Một số ví dụ cho phương pháp này như Tapia and Rojas Đã đề xuất một phương pháp nhận diện dựa trên cây trải dài²⁴ và ký tự chủ đạo²⁵, Zanibbi cho ra đời phương pháp nhận diện bằng một chuỗi các bước biến đổi cây[7]. Phương pháp này khá hoàn thiện nhưng vẫn có một số vướng mắc như việc nhận diện ký tự một cách phi ngữ cảnh vẫn mang đến sự thiếu tự nhiên, hay phương pháp đọc từ trái sang phải vẫn để lại nhiều lỗ hổng trong việc nhận diện biểu thức.
- Nhận diện dựa trên ngữ pháp toán học. Ở hướng này, ngữ pháp được đưa vào quá trình nhận diện, ví dụ như sử dụng ngữ pháp để hậu xử lý, loại bỏ, hiệu chỉnh các ý tự nhận diện sai hoặc sử dụng các mô hình ngữ pháp để dự đoán ký tự tiếp theo trong biểu thức.

²¹Thuật ngữ tiếng Anh: Optical Character Recognition, viết tắt OCR.

²²Thuật ngữ tiếng Anh: Mathematical Expression, viết tắt ME.

²³Thuật ngữ tiếng Anh: Segmentation.

²⁴Thuật ngữ tiếng Anh: Spanning tree.

²⁵Thuật ngữ tiếng Anh: Dominate symbol.

2 WHOOOOOO : Watch, Attend and Parse: An End-to-end Neural Network Based Approach to Handwritten Mathematical Expression Recognition

Khác với những bài báo đi trước sử dụng phương pháp phân vùng hay dựa trên ngữ pháp, bài báo này sử dụng phương pháp mang tên Watch, Attend, Parse dựa trên bài toán chú thích cho ảnh²⁶. Hệ thống này là sự kết hợp giữa CNN, RNN và một hệ thống ANN (Cụ thể là GRU²⁷), CNN đặc trưng cho Watcher sẽ trích đặc trưng ảnh, ANN đặc trưng cho Attend mang nhiệm vụ điều hướng cho mạng biết được vị trí nào sẽ tập trung vào và RNN đặc trưng cho Parser mang nhiệm vụ sinh ra chuỗi Latex chính là đầu ra của hệ thống

Cấu trúc mạng này được chia thành ba phần:

2.1 Watcher

Bộ phận này có bản chất là một hệ thống mạng nơ-ron tích chập đầy đủ (FCN)²⁸ có cấu tạo là các lớp tích chập và các lớp pooling xếp chồng lên nhau. Bộ phận này nhận đầu vào là ảnh cần nhận diện và đầu ra là các vector đặc trưng ứng với mỗi pixel của ảnh.

2.2 Attend

Bộ phận này hoạt động giống như một hệ thống tổng hợp thông tin, nó có chức năng lấy dữ liệu ký tự vừa được dự đoán từ Parser (sẽ được giải thích bên dưới), từ các vector đặc trưng đã có được từ Watcher và ghi nhận các vị trí đã được xử lý trong ảnh, từ đó dự đoán vị trí tiếp theo để xử lý.

2.3 Parser

Bản chất của bộ phận này là một mạng GRU nhận dữ liệu đầu vào từ hai bộ phận còn lại, từ đó sinh ra chuỗi Latex

3 Wenhao He[1]: Context-aware Recognition

Một quá trình nhận dạng biểu thức toán học chuẩn bao gồm 2 giai đoạn: giai đoạn đầu tiên là phân tách và nhận dạng ký hiệu²⁹, giai đoạn thứ hai là phân tích cấu trúc³⁰. Nhiều công trình nghiên cứu trước đây về nhận dạng biểu thức toán học cũng thực hiện theo phương pháp như vậy. Nhược điểm của những phương pháp dạng này là thực hiện hai quá

²⁶Thuật ngữ tiếng anh: Image Captioning

²⁷Viết tắt: Gated recurrent unit

²⁸Thuật ngữ tiếng anh: Fully Convolution Network

²⁹symbol segmentation and recognition

³⁰structure analysis

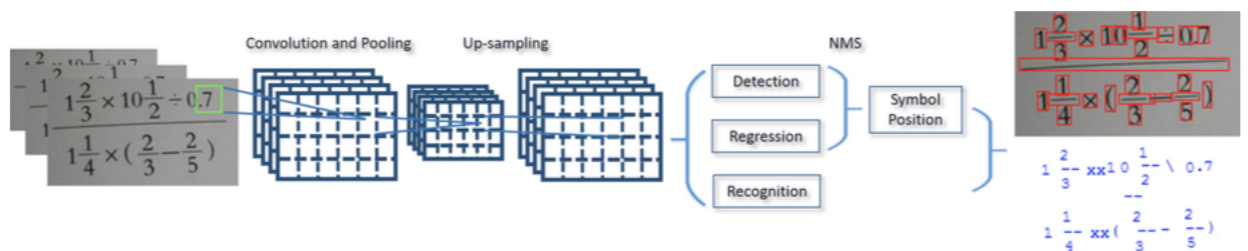
trình trên một cách độc lập, do đó thông tin cấu trúc³¹ của biểu thức không được đưa vào quá trình nhận dạng. Như vậy sẽ dẫn đến lỗi[1].

WenHao He cùng các cộng sự của ông đã đề xuất một phương pháp dựa trên CNN, kết hợp với cách học đa nhiệm vụ³², cố gắng kết hợp hai giai đoạn chuẩn của quá trình nhận dạng biểu thức toán lại với nhau.

Phương pháp mà các ông đưa ra đảm bảo thông tin cấu trúc của biểu thức được đưa vào quá trình nhận dạng và thể hiện trong các ma trận đặc trưng³³[4].

Đối với các phương pháp trước, biểu thức phải được phân tách thành những ký hiệu rồi từng ký hiệu này mới được đưa qua bộ nhận dạng. Một số phương pháp phân tách ký hiệu thường được sử dụng như: phân tích thành phần liên thông³⁴, cắt dựa trên các phép chiếu³⁵[8]. Tuy nhiên với phương pháp mới này, cả ảnh của biểu thức toán học được đưa qua bộ nhận dạng, chính vì vậy mà thông tin cấu trúc của biểu thức được bảo toàn.

Dưới đây là mô hình học của phương pháp này:



Hình 2: Mô hình học được đề xuất [1].

Ảnh đầu vào sẽ qua một số lớp tích chập, down-sampling và up-sampling để tạo ra một feature map. Feature map này sẽ được gửi đến ba nhiệm vụ, mỗi nhiệm vụ sẽ tạo ra 1 feature map có cùng kích thước với feature map đầu vào.

Giả sử một điểm i được cho đặt tại toạ độ (w_i, h_i) của feature map được tạo ra bởi các nhiệm vụ.

- **Nhiệm vụ phát hiện** (Detection task) sẽ cho ra một con số s thể hiện độ tin cậy rằng một ký hiệu được đặt tại i .
- **Nhiệm vụ hồi quy** (Regression task) cho ra một vector 4 chiều x_1, y_1, x_2, y_2 thể hiện thông tin về bounding box của ký hiệu được đặt tại i .
- **Nhiệm vụ nhận dạng** (Recognition task) gán nhãn cho ký hiệu đặt tại i cùng với xác suất của nhãn đó.

³¹structure information hay context information

³²multi-task learning

³³feature map

³⁴connected components

³⁵projection cutting

Như vậy nhiệm vụ phát hiện và hồi quy được thiết kế để định vị trí của ký hiệu trong biểu thức toán học, nhiệm vụ nhận dạng quyết định xem đó là ký hiệu gì.

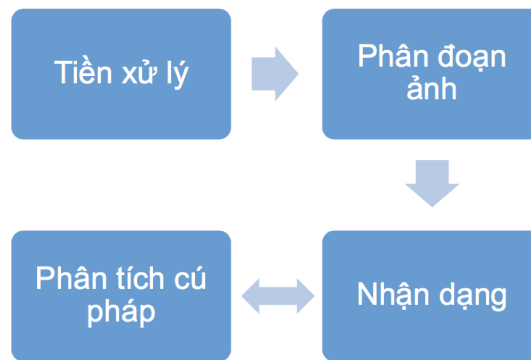
Phương pháp nhận dạng như trên có thể giải quyết cả những trường hợp là thách thức đối với các phương pháp phân tách và nhận dạng ký hiệu trước đây, cụ thể đó là vấn đề gom nhóm ký hiệu đối với các ký hiệu nhiều phần nhỏ, tách những ký hiệu bị dính nhau.

4 QAK[2]

QAK là tên dự án được thực hiện bởi nhóm sinh viên khoá 2011 cũng về đề tài nhận dạng biểu thức toán học.

Phương pháp nhóm sinh viên này xây dựng dựa trên nghiên cứu chính từ hai công trình của Aderson[9] và Zanibbi[7]. Do đó, qua trình nhận dạng biểu thức toán học vẫn đi theo 2 bước chuẩn đó là phân tách ký hiệu và phân tích cấu trúc như đã trình bày ở mục 3.

Dưới đây là mô hình phương pháp mà nhóm đã đề xuất:



Hình 3: Quy trình nhận dạng.

- Trong bước tiền xử lý, nhóm áp dụng một số kỹ thuật trong xử lý ảnh như: loại nhiễu, ảnh nhị phân,... để tăng cường chất lượng ảnh, hỗ trợ cho bước phân đoạn.
- Trong bước phân đoạn ảnh, nhóm dùng kỹ thuật chính là cắt theo hình chiếu[8] và phân tích thành phần liên thông[8] cho toán tử lấy căn. Mục tiêu của giai đoạn này là phân tách ảnh chứa biểu thức toán học ban đầu ra thành các mảnh ảnh, mỗi mảnh chỉ chứa 1 ký hiệu toán học. Ngoài ra, ở bước này một cấu trúc cây được xây dựng lưu thông tin của các mảnh ảnh, hỗ trợ cho quá trình phân tích ngữ pháp sau này.
- Ở bước nhận dạng, nhóm sinh viên khoá 2011 đã sử dụng một kiến trúc mạng tương tự Lenet-5[4].
- Với phân tích cú pháp, nhóm sử dụng tập luật **văn phạm phi ngữ cảnh**³⁶ do chính nhóm đề xuất để giới hạn kết quả đầu ra của quá trình nhận dạng, từ đó tăng khả năng nhận dạng đúng biểu thức.

³⁶Thuật ngữ tiếng Anh: Context-free grammar

Chương 4 Mô hình đề xuất

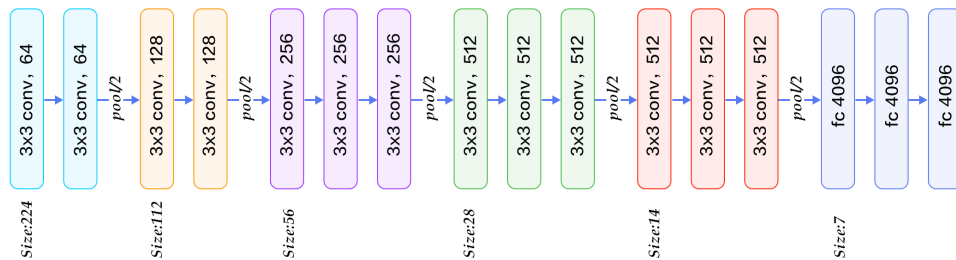
1 Tổng quan

Nhóm quyết định sử dụng mạng SSD³⁷ để phát hiện các ký tự trong ảnh và sau đó sử dụng bộ phân tích cú pháp DRACULAE để sinh ra chuỗi LaTeX.

2 Phát hiện ký tự

2.1 Mạng cơ sở

Về cấu trúc của SSD, nhóm đã sử dụng mạng VGG16 [trích dẫn] cho phần mạng cơ sở. Cấu trúc của mạng VGG16 gồm tổ hợp các lớp tích chập và lớp pooling xếp chồng lên nhau, ở cuối mạng có các lớp liên kết đầy đủ³⁸ để giảm kích thước tensor và cuối cùng là xuất ra vector có số chiều phù hợp (có số chiều bằng với số lớp đối tượng cần dự đoán).



Hình 4: Cấu tạo mạng VGG16

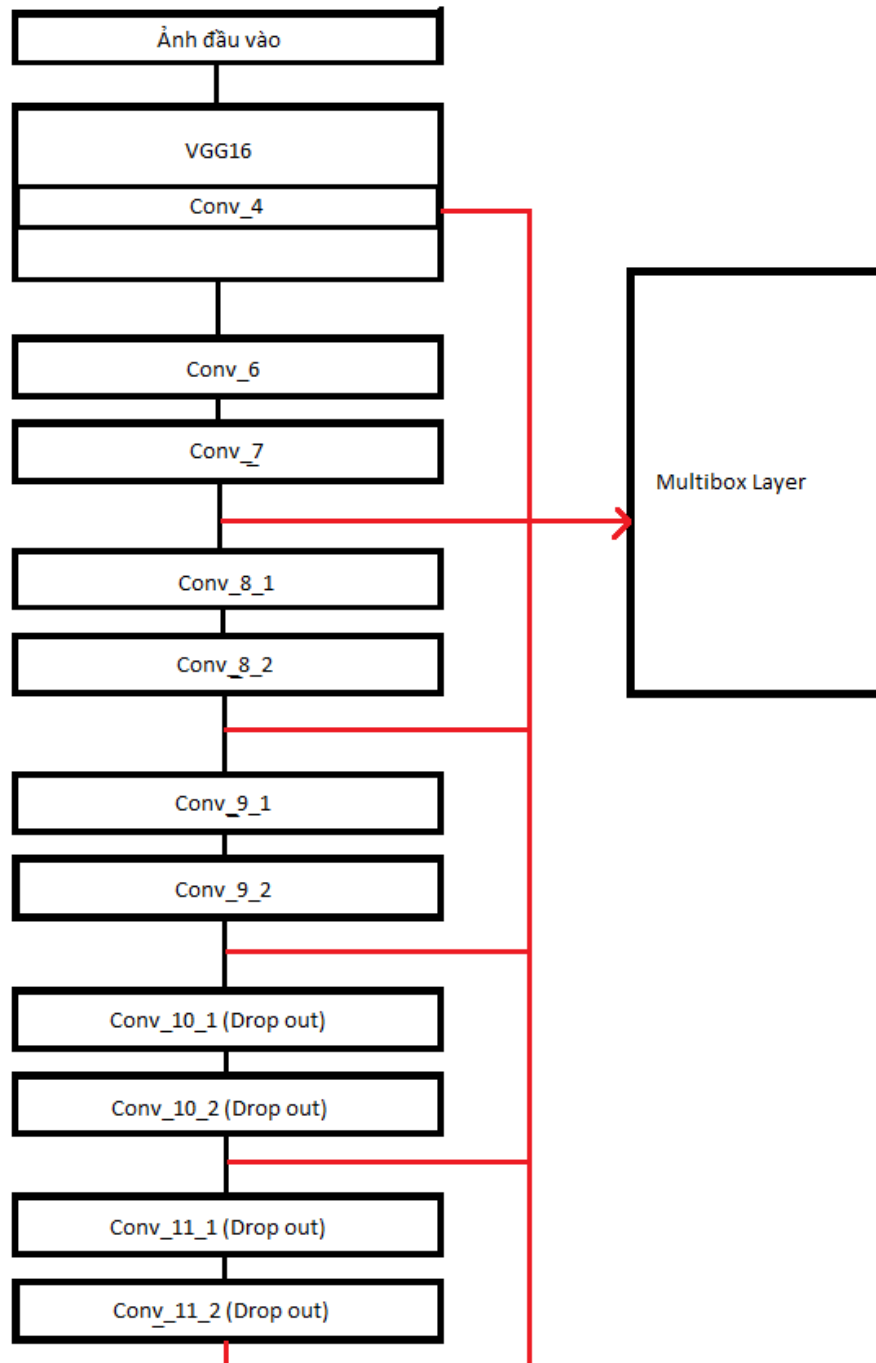
Như hình vẽ bên trên, các lớp tích chập và lớp pooling được gom lại thành các cụm, sau hai đến ba lớp tích chập là một lớp pooling. Số kênh của tensor chạy trong mạng được nâng dần từ 64 lên 512 về phía cuối mạng. Khi đưa mạng cơ sở vào mạng SSD, các lớp liên kết đầy đủ sẽ được bỏ đi và bộ phận phân loại sẽ được chèn vào vị trí tương ứng. Khi đó, ta cần phải chuyển model của mạng VGG sang mạng SSD, phần này sẽ được đề cập chi tiết hơn trong phần hiện thực.

2.2 Thân mạng SSD

Khi tiếp hợp vào mạng SSD, các lớp liên kết đầy đủ của VGG16 được lược bỏ, thay vào đó là những lớp tích chập, mô hình cụ thể của mạng SSD nhóm đề xuất được thể hiện ở ảnh dưới.

³⁷Thuật ngữ tiếng Anh: Single Shot Detector

³⁸Thuật ngữ tiếng Anh: Fully connected layer



Hình 5: Sơ đồ các lớp trong mạng SSD

Trong đó Khối VGG16 là phần mạng VGG16 đã được lược bỏ các lớp liên kết đầy đủ. Các khối còn lại được thể hiện trong bảng sau:

Khối	Số kênh đầu vào	Số kênh đầu ra	Kích thước nhân	Chèn thêm (Padding)	Bước (Stride)
Conv_6	512	1024	3	1	1
Conv_7	1024	1024	3	6	1
Conv_8_1	1024	256	1	0	1
Conv_8_2	256	512	3	1	2
Conv_9_1	512	128	1	0	1
Conv_9_2	128	256	3	1	2
Conv_10_1	256	128	1	0	1
Conv_10_2	128	256	3	1	2
Conv_11_1	256	128	1	0	1
Conv_11_2	128	256	3	0	1

Riêng khối Conv_4 trong ảnh là lớp tích chập thứ 10 trong mạng VGG.

Các lớp tích chập đều sử dụng hàm truyền là hàm ReLu, riêng ở sau các lớp Conv_10_1, Conv_10_2, Conv_11_1, Conv_11_2 thì có thêm một lớp dropout với hệ số dropout là 0.2.

Các feature map được trích ra để thực hiện phát hiện, phân loại từ các mẫu tên màu đỏ trên hình, cụ thể:

- Ngay sau khối tích chập - pooling 512 thứ nhất trong mạng VGG (hay sau lớp pooling của lớp tích chập thứ 10 của mạng VGG - hay Conv_4).
- Sau các khối tích chập Conv_7, Conv_8, Conv_9, Conv_10, Conv_11.

Đối với mạng SSD300[trích dẫn] nguyên thủy, kích thước của các feature map ứng với số dự đoán được thể hiện ở bảng sau:

Feature map Sinh bởi	Kích thước	Tỉ lệ diện mạo	Số dự đoán
Conv_4	$512 \times 38 \times 38$	$\{1, \frac{1}{2}, 2\}$	5776
Conv_7	$1024 \times 19 \times 19$	$\{1, \frac{1}{2}, 2, \frac{1}{3}, 3\}$	2166
Conv_8	$512 \times 10 \times 10$	$\{1, \frac{1}{2}, 2, \frac{1}{3}, 3\}$	600
Conv_9	$256 \times 5 \times 5$	$\{1, \frac{1}{2}, 2, \frac{1}{3}, 3\}$	150
Conv_10	$256 \times 3 \times 3$	$\{1, \frac{1}{2}, 2\}$	36
Conv_11	$256 \times 1 \times 1$	$\{1, \frac{1}{2}, 2\}$	4

Như vậy, ứng với mỗi ảnh, hệ thống sinh ra 8732 dự đoán. Tương ứng, bộ phận mã hóa của mạng SSD300 cũng có những thông số:

- Danh sách kích thước feature map: (38, 19, 10, 5, 3, 1)
- Danh sách kích thước các ô chuẩn (với $s_{min} = 0.2$ và $s_{max} = 0.9$): (30, 60, 111, 162, 213, 264, 315)
- Ngưỡng Jaccard: 0.5.

2.3 Các thay đổi

Do bản chất mạng SSD300 gặp rất nhiều khó khăn trong việc nhận dạng ký tự nhỏ do một số lý do:

- Trong quá trình kết hợp, các ô bọc "đáp án" kết hợp với ô chuẩn phải có chỉ số Jaccard lớn hơn 0.5, nhưng kích thước tối thiểu của ô chuẩn lại là 30×30 pixel, vì vậy nhiều ký tự nhỏ không thể kết hợp được (Trường hợp này xảy ra với hầu hết những ảnh có chữ nhỏ, ví dụ như ảnh dưới)

$$p = 2 \times l - q$$

Hình 6: Ảnh ví dụ chữ nhỏ (với kích thước thật)

Khi không thể kết hợp được thì vùng ký tự đó được hiểu là phần nền, điều này gây bất lợi lớn cho quá trình huấn luyện do vùng ảnh không trống nhưng lúc là nền lúc lại được gắn nhãn.

- Nhiều ký tự có tỉ lệ diện mạo đặc biệt (Ví dụ như ký tự dấu $-$ trong phân số hoặc \int) cũng gặp nhiều khó khăn trong việc kết hợp do không đạt được ngưỡng Jaccard.

Vì vậy, nhóm đã sửa một số thông số để tìm cách giải quyết vấn đề trên, trong nhiều bộ thông số thì có 3 mô hình nổi bật.

Giảm kích thước các ô chuẩn

Mô hình này dựa trên mô hình SSD300 nguyên thủy, chỉ chỉnh sửa kích thước các ô bọc thông qua chỉnh hai thông số $s_{min} = 0.08$ và $s_{max} = 0.5$, từ đó danh sách kích thước các ô chuẩn được thay đổi thành (9, 24, 54, 84, 114, 144, 174). So với mô hình nguyên thủy, mô hình này không thể kết hợp tốt các ký tự lớn (lớn hơn khoảng 200×200 pixel). Tuy vậy, đối với bài toán về phát hiện các ký tự trong một ảnh, thì kích thước các ký tự không thể quá lớn, nên những ký tự đặc biệt lớn như vậy có thể được xem xét ở độ ưu tiên thấp hơn.

Việc thay đổi tham số ở mô hình này không làm ảnh hưởng lớn đến cấu trúc mạng cũng như quá trình mã hóa, giải mã. Số lượng dự đoán không thay đổi.

Mục đích chính của thay đổi này là để quá trình kết hợp diễn ra tốt hơn khi kích thước của ô chuẩn nhỏ nhất là 9×9 , điều này giúp cho những ký tự rất nhỏ vẫn có thể được kết hợp, tránh hiện tượng bỏ sót ký tự như mạng SSD300 nguyên thủy.

Giảm kích thước các ô chuẩn và tăng kích thước ảnh đầu vào

Để có được mô hình này, nhóm cần phải có hai thay đổi so với mạng SSD300 nguyên thủy:

- Thay vì sử dụng ảnh 300×300 như nguyên thủy, ảnh sẽ được phóng to lên kích thước 500×500 .
- Như ở mô hình đề xuất trước, nhóm thay đổi $s_{min} = 0.08$ và $s_{max} = 0.5$ từ đó thu được danh sách kích thước (15.0, 40, 80, 120, 160., 200, 240, 280).

Những thay đổi trên có làm thay đổi đến cấu trúc mạng SSD, khi tăng kích thước ảnh đầu vào nghĩa là kích thước của các feature map thu thập được cũng tăng lên, từ đó kéo số lượng dự đoán trong một ảnh tăng lên như bảng dưới.

Feature map Sinh bởi	Kích thước	Tỉ lệ diện mạo	Số dự đoán
Conv_4	$512 \times 63 \times 63$	$\{1, \frac{1}{2}, 2\}$	15876
Conv_7	$1024 \times 32 \times 32$	$\{1, \frac{1}{2}, 2, \frac{1}{3}, 3\}$	6144
Conv_8	$512 \times 16 \times 16$	$\{1, \frac{1}{2}, 2, \frac{1}{3}, 3\}$	1536
Conv_9	$256 \times 8 \times 8$	$\{1, \frac{1}{2}, 2, \frac{1}{3}, 3\}$	384
Conv_10	$256 \times 4 \times 4$	$\{1, \frac{1}{2}, 2\}$	144
Conv_11	$256 \times 2 \times 2$	$\{1, \frac{1}{2}, 2\}$	64

Như vậy với mỗi ảnh, hệ thống sinh ra 24148 dự đoán, gấp gần 3 lần so với ảnh SSD300 nguyên thủy. Với những thay đổi này ô boc thậm chí còn có thể kết hợp với những ô chuẩn nhỏ hơn, bên cạnh đó, mật độ các ô chuẩn cũng dày đặc hơn nhiều so với mạng SSD nguyên thủy.

Giảm kích thước các ô chuẩn, tăng kích thước ảnh đầu vào và thêm lớp tích chập

Mô hình này thừa kế các thay đổi ở hai mô hình trước, tuy nhiên ở phần thân mạng SSD, nhóm đã gắn thêm 2 lớp tích chập ở cuối với các tham số:

Khối	Số kênh đầu vào	Số kênh đầu ra	Kích thước nhân	Chèn thêm (Padding)	Bước (Stride)
Conv_12_1	256	128	1	0	1
Conv_12_2	128	256	3	0	1

Hai lớp tích chập này cũng có hàm truyền là hàm relu và đi chung với lớp drop out với hệ số drop out là 0.2.

Thay đổi này cũng làm thay đổi cấu trúc mạng:

Feature map Sinh bởi	Kích thước	Tỉ lệ diện mạo	Số dự đoán
Conv_4	$512 \times 63 \times 63$	$\{1, \frac{1}{2}, 2\}$	15876
Conv_7	$1024 \times 32 \times 32$	$\{1, \frac{1}{2}, 2, \frac{1}{3}, 3\}$	6144
Conv_8	$512 \times 16 \times 16$	$\{1, \frac{1}{2}, 2, \frac{1}{3}, 3\}$	1536
Conv_9	$256 \times 8 \times 8$	$\{1, \frac{1}{2}, 2, \frac{1}{3}, 3\}$	384
Conv_10	$256 \times 4 \times 4$	$\{1, \frac{1}{2}, 2\}$	160
Conv_11	$256 \times 2 \times 2$	$\{1, \frac{1}{2}, 2\}$	40
Conv_12	$256 \times 2 \times 1$	$\{1, \frac{1}{2}, 2\}$	10

Lớp Conv_12 được thêm vào cuối mạng với mục đích tăng độ dài danh sách feature map đưa vào lớp Multibox, điều này giúp trong quá trình kết hợp, các ký tự lớn hơn có thể được kết hợp.

3 Phân tích cú pháp

Bộ phân tích cú pháp này có nhiệm vụ chuyển dữ liệu dự đoán là các ô bọc sang dữ liệu biểu thức toán học dạng latex. Trong quá trình phân tích cú pháp, biểu thức được lưu trữ dưới dạng cây BST³⁹ là một cấu trúc cây dựa trên các đường cơ sở, các nhánh của cây thể hiện một phân vùng bên trên, bên dưới, ... của nút và là một phân vùng ảnh có đường cơ sở riêng và Lexed - BST là một cấu trúc cây tương tự như BST nhưng thay vì chỉ biểu thị vị trí bên trên, bên dưới thì cây biểu thị quan hệ phân số, chỉ số trên, chỉ số dưới, ...

3.1 Sinh cây BST

Đây là giai đoạn phức tạp nhất và cũng là quan trọng nhất, ở giai đoạn này, dữ liệu đầu ra từ mạng SSD là danh sách các ô bọc sẽ được phân tích và xây dựng một cây BST. Công việc này có thể được phân ra thành bốn công đoạn: Tìm ký tự chủ đạo, xác định những ký tự trên đường cơ sở, tái phân vùng và xử lý nút con

Trước khi tiến hành sinh cây, ta cần phải phân loại các ký tự vào những thể loại khác nhau, điều này có mục đích giúp xác định vị trí trọng tâm tốt hơn, xác định được vùng ảnh gắn với chỉ số trên, chỉ số dưới và giúp phân biệt một số loại ký tự có thể có cấu trúc đặc biệt trong biểu thức (ví dụ như ký tự \sum thì có thể có ký tự nằm trên và bên dưới, nhưng ký tự *sin* thì gần như không bao giờ có). Các ký tự được phân vào các lớp như bảng dưới:

³⁹Viết tắt: baseline Syntax Tree: Cây cú pháp dựa trên đường cơ sở



Lớp ký tự	Tung độ trọng tâm	Bên dưới	Bên trên	Chỉ số dưới	Chỉ số trên
Non-Script Các ký tự phép tính (+, -, →, ...)	$\frac{1}{2}H$	$\frac{1}{2}H$	$\frac{1}{2}H$	-	-
Open Bracket Ký tự mở ngoặc (cH	$\min(Y)$	$\max(Y)$	-	-
Root Ký tự căn $\sqrt{\quad}$	cH	$\min(Y)$	$\max(Y)$	tH	$H - tH$
Variable Range Ký tự có thể viết ở trên hoặc dưới $\sum, \int, \lim \dots$	$\frac{1}{2}H$	tH	$H - tH$	tH	$H - tH$
Ascender Ký tự nổi lên A..Z, 0..9, b,d,f ...	cH	tH	$H - tH$	tH	$H - tH$
Descender Ký tự chìm xuống g, p, y, j, ..., b,d,f ...	$H - cH$	$\frac{1}{2}H + t\frac{1}{2}H$	$H - t\frac{1}{2}H$	$\frac{1}{2}H + t\frac{1}{2}H$	$H - t\frac{1}{2}H$
Centered Ký tự trung tâm o, n, u, },) ...	$\frac{1}{2}H$	tH	$H - tH$	tH	$H - tH$

Trong đó:

- Các cột thứ hai trở đi là tọa độ của một số điểm, ngưỡng đặc biệt tính từ góc trái bên dưới của ký tự. Hệ trục tọa độ có trục tung hướng lên và trục hoành hướng sang phải. Cột "Tung độ trọng tâm" thể hiện tung độ y của trọng tâm của ký tự, cột bên dưới và bên trên thể hiện ngưỡng trên và ngưỡng dưới của ký tự, các ký tự vượt ngoài ngưỡng đó sẽ được xem là nằm bên trên hoặc nằm bên dưới ký tự đang xét (Ví dụ như một ký tự nào đó có tung độ trọng tâm lớn hơn $\max(Y)$ của một ký tự thuộc lớp *Root* thì sẽ được xem là nằm bên trên). Cột chỉ số trên, chỉ số dưới thể hiện ngưỡng để được phân vào vùng chỉ số trên, chỉ số dưới. Riêng hai lớp Non-script và Open Bracket có giá trị hai cột này không xác định là do những ký tự của lớp này không có chỉ số trên, dưới.
- H là chiều cao của ký tự.
- c (Viết tắt của centroid ratio): là tham số để xác định trọng tâm cho một số ký tự đặc biệt (ví dụ như các ký tự ngoặc hoặc chìm xuống dưới), tùy người viết chữ có thể có những giá trị phù hợp khác nhau.

- t Là ngưỡng cho chỉ số trên, dưới và chỉ phần bên trên, bên dưới. Tham số này biểu thị sự nhạy cảm đối với những ký tự trên dưới, tham số này càng lớn thì các ký tự khi được phân tích vị trí càng dễ được xếp vào chỉ số trên, chỉ số dưới.

Bảng trên thể hiện quy tắc để kiểm tra hai ký tự có vị trí tương đối như thế nào với nhau, cụ thể:

- B chứa trong A (dành cho trường hợp ký tự căn: $\sqrt{}$) khi B có tung độ nằm trong khoảng từ ngưỡng chỉ số dưới đến ngưỡng chỉ số trên của A và hoành độ trọng tâm của B nằm trong khoảng $(Min(X_B), Max(X_B))$.
- A và B liên kề nếu như trọng tâm của B có tung độ nằm trong khoảng từ ngưỡng chỉ số dưới đến ngưỡng chỉ số trên của A và hoành độ trọng tâm của B lớn hơn $Max(X_A)$.
- B là chỉ số trên của A nếu như trọng tâm của B lớn hơn ngưỡng chỉ số trên của A.
- B là ký tự phân số thống trị A khi trọng tâm của A có hoành độ nằm trong khoảng $(Min(X_B), Max(X_B))$

Trước khi tiến hành sinh cây BST, danh sách ô bọc cần phải được sắp xếp theo chiều tăng dần hoành độ của mép phải ô bọc x (hay cụ thể hơn là $min(x)$).

Tìm ký tự chủ đạo

Trong hầu hết trường hợp, ký tự chủ đạo là ký tự nằm bên trái nhất, tuy vậy, ta vẫn cần phải kháng được những trường hợp người dùng viết chữ không chuẩn gây thụt ra ngoài, hoặc những ký tự có dạng như:

$$\sum_{n=100000} a$$

Thì ký tự chủ đạo là \sum nhưng ký tự bên trái nhất lại là n , Draculae được thế kế để có thể chống lại những trường hợp như vậy.

Để tìm ký tự chủ đạo, ta cần thực hiện những công đoạn sau:

- Trước hết, ta tìm ký tự nằm bên trái nhất dựa trên $min(X)$ của ô bọc, xem tất cả các ký tự là nút con của nút gốc.
- Kiểm tra ký tự hiện tại có phải là Variable Range (lớp này có đặc tính có thể viết cả ở bên trên và bên dưới, đây là nguồn gốc gây khó khăn trong việc xác định ký tự chủ đạo). Nếu đúng thì ta kiểm tra xem ký tự đó có bị thống trị bởi một phân số hay không, nếu có thì ký tự chủ đạo là ký tự phân số, nếu không thì ký tự chủ đạo là ký tự Variable Range vừa tìm được. Nếu ký tự vừa tìm được không phải là Variable Range, thì ta tìm kiếm ký tự Variable Range nằm bên trái nhất và sang bước tiếp theo.



- Nếu không còn ký tự Variable Range nào khác, thì ký tự bên trái nhất vừa tìm được ở bước trên được kiểm tra bị thống trị bởi phân số hay không và trả về ký tự chủ đạo (ký tự đó hoặc ký tự phân số). Nếu còn một ký tự Variable Range khác trong biểu thức, ta tiếp tục với bước kế tiếp.
- Ta kiểm tra xem ký tự Variable Range khác đó có ký tự nào nằm liền kề bên trái không. Nếu không thì ký tự Variable Range được kiểm tra xem có bị thống trị bởi ký tự phân số hay không và chọn ký tự chủ đạo ký tự phù hợp. Nếu tồn tại ký tự liền kề bên trái, thì ký tự bên trái nhất được kiểm tra có bị thống trị bởi phân số không và chọn ký tự chủ đạo phù hợp.

Xác định ký tự trên đường cơ sở

Sau khi đã có được ký tự chủ đạo, ta cần phải tìm những ký tự liền kề theo chiều ngang với ký tự vừa tìm được, những ký tự này chính là ký tự trên đường cơ sở. Để xác định các ký tự trên đường cơ sở, ta lấy ký tự chủ đạo vừa xác định được để xét:

1. Nếu như ký tự đang xét là ký tự cuối cùng hoặc ký tự duy nhất được xét, thì chuyển sang bước kế tiếp (bước tái phân vùng).
2. Ta phân vùng các ký tự nằm ở vùng bên trên, bên dưới, góc trên bên trái và góc dưới bên trái (không phân vùng các ký tự nằm bên phải) vào nút con của nút đang xét.
3. Nếu như ký tự đang xét thuộc lớp Non-script, ký tự đó được đánh dấu nằm trên đường chủ đạo, ta quay lại tìm ký tự chủ đạo trong các ký tự còn lại và xác định các ký tự trên đường cơ sở với ký tự chủ đạo đó (quay lại bước 1).
4. Ta kiểm tra danh sách các ký tự còn lại xem có ký tự nào liền kề bên phải với ký tự đang xét hay không, nếu có, ta kiểm tra ký tự đó có bị thống trị bởi ký tự khác hay không và quay lại bước 1 với ký tự được xét là ký tự thống trị nhất (hoặc chính ký tự liền kề nhất vừa tìm được đối với trường hợp nó không bị thống trị bởi ký tự nào).
5. Nếu hệ thống chạy đến bước này thì có nghĩa là trong các ký tự còn lại, không còn ký tự nào liền kề với ký tự đang xét, vì vậy ta tiến hành phân vùng đưa các ký tự còn lại vào nút con góc trên bên phải (chỉ số trên) và góc dưới bên phải (chỉ số dưới) của nút đang xét.

Sau khi thực hiện bước này lần đầu tiên, ta thu được một cây có ba tầng, tầng đầu chứa nút gốc, tầng thứ hai chứa các ký tự nằm trên đường cơ sở, tầng thứ ba chứa các nút con của các nút chứa ký tự nằm trên đường cơ sở, các nút con này có thể chứa một hoặc nhiều ký tự, chưa là một cây.

Tái phân vùng

Ở bước trước, nút con được phân vào vùng góc trên bên trái và góc dưới bên trái, điều này không có ý nghĩa đối với ngữ nghĩa của biểu thức toán học và khó khăn trong việc sinh cây Lexed BST sau này. Vì vậy ở bước này, các nút con sẽ được thay đổi vị trí cho phù hợp với ngữ nghĩa của biểu thức, chương trình sẽ kiểm tra từng cặp ký tự liền kề theo thứ tự từ trái sang phải, ta gọi ký tự bên trái là A và ký tự bên phải là B :

- Nếu B không thuộc lớp Variable Range, và B không có nút con ở vùng bên trên, ta gắn nút góc trên bên trái của ký tự B vào nút góc trên bên phải (chỉ số trên) của A .
- Nếu B thuộc lớp Variable Range thì ta sẽ phân vùng những ký tự liền kề với ký tự đầu tiên trong nút con góc trên bên trái của B vào vùng bên trên của B .
- Nếu A thuộc lớp Variable Range thì ta phải tiến hành gộp các vùng góc trên bên trái, bên trên và góc trên bên phải thành một vùng bên trên duy nhất.

Sau đó ta tiến hành làm tương tự với vùng bên dưới. Với cách phân vùng này, với một số biểu thức đặc biệt hoặc nhập nhầm thì việc phân vùng sẽ không chính xác, để khắc phục điều này, ta cần phải có thêm công đoạn hậu xử lý.

Đệ quy

Tới bước này, ta đã có một cây có ba tầng với tầng thứ hai có phân vùng nút con khá chuẩn xác. Nhiệm vụ bây giờ là phải tiếp tục phân vùng cho các nút con của các ký tự trên đường cơ sở. Việc này được thực hiện bằng cách xem từng nút con là từng vùng ảnh mới riêng biệt và thực hiện lại ba bước trên để thu được các cây con.

Sau quá trình này, ta thu được một cây BST biểu thị vị trí tương đối giữa các nút. Cây này tiếp tục được đi qua [Lexical Pass] để thu được Lexed - BST.

3.2 Sinh cây Lexed - BST

Sau khi đã thu được cây BST từ danh sách các ô bọc, nhiệm vụ còn lại là không quá phức tạp. Trước khi tạo được chuỗi Latex, ta cần phải sinh ra được một cây Lexed - BST, cây này chứa các thông tin như "nhóm ký tự" và cấu trúc biểu thức:

- Nhóm ký tự: Ban đầu, các ký tự được gom thành từng nhóm có nghĩa (Ví dụ các ký tự "s", "i", "n" thì sẽ được gộp lại thành một ký tự "sin" duy nhất).
- Ta sử dụng các luật khác nhau để sinh ra cây Lexed - BST.



Luật sinh

Để chuyển từ cây BST sang cây lexed - BST, nhóm có đề xuất một số luật sinh:
Gọi ký tự đang xét là A

- Nếu cây BST B nằm trong nút con góc phải bên dưới của A thì sinh ra một nút $\text{sub}(A, B)$
- Nếu cây BST B nằm trong nút con góc phải bên trên của A thì sinh ra một nút $\text{sup}(A, B)$
- Nếu cây BST B nằm trong nút con bên trong của A và A là ký tự căn ($\sqrt{}$) thì sinh ra một nút $\text{sqrt}(B)$
- Nếu A có nút con bên trên hoặc bên dưới (hoặc cả hai), thì tùy thuộc vào ký tự A mà ta sinh ra các nút khác nhau (Ví dụ như các ký tự \sum , \lim , \prod , \int , \rightarrow và quan trọng nhất là dấu gạch ngang trong phân số)

Sau khi sinh được cây Lexed - BST, hệ thống sẽ duyệt qua cây theo chiều tiền thứ tự để sinh mã Latex

Chương 5 Hiện thực, đánh giá

Để xây dựng hệ thống nhận diện biểu thức toán học viết tay, nhóm cần hoàn thành ba công việc chính bao gồm thu thập dữ liệu, huấn luyện mạng SSD và xây dựng chương trình sinh mã latex.

1 Chuẩn bị dữ liệu

2 Hiện thực hệ thống

2.1 Quá trình tiền huấn luyện

Như đã trình bày ở trên, mạng SSD mà nhóm lựa chọn có sử dụng mạng cơ sở là VGG16, đây là một mạng phân lớp⁴⁰. Để đạt được kết quả huấn luyện cho mạng SSD tốt hơn, tác giả đã khuyến cáo sử dụng phương pháp học chuyển tiếp⁴¹

Học chuyển tiếp

Để thực hiện phương pháp này, ta cần phải chuẩn bị một tập dữ liệu riêng gồm các ảnh kích thước 32×32 . Các ảnh này được lấy từ tập huấn luyện của đề tài QAK.

⁴⁰Thuật ngữ tiếng anh: Classification

⁴¹Thuật ngữ tiếng anh: Transfer Learning

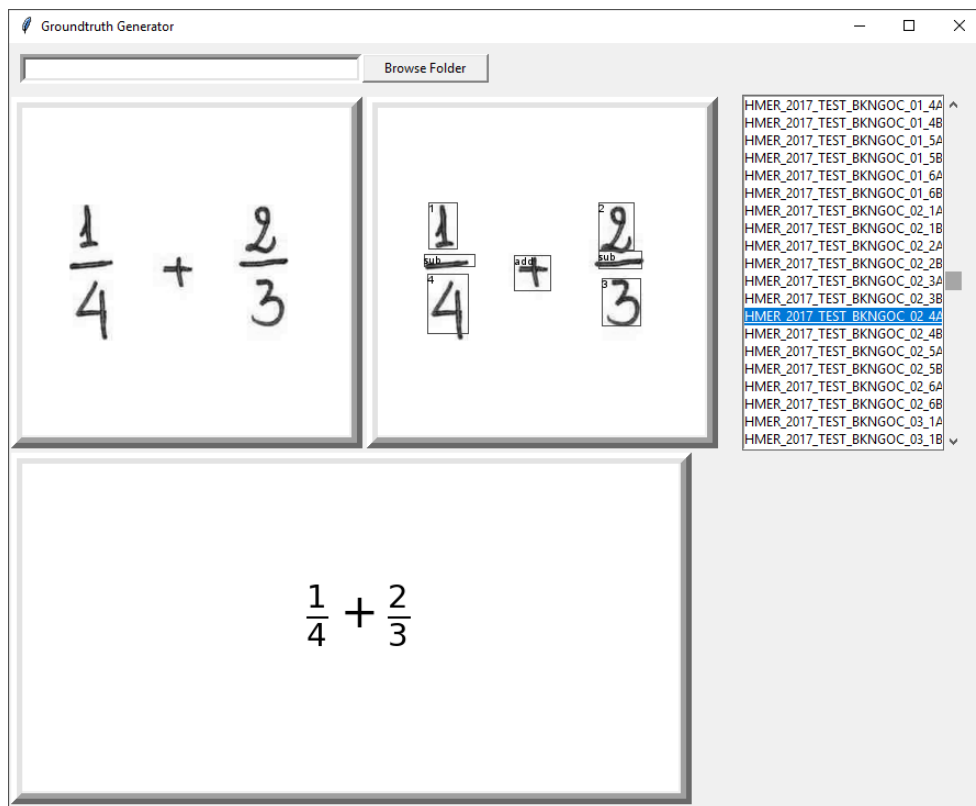
Khi đã thu được một mô hình phù hợp (Huấn luyện tập tiền huấn luyện đã hội tụ), nhiệm vụ kế tiếp là ta phải chuyển các tham số đã học được sang một mô hình của mạng SSD.

Quá trình này không quá phức tạp, nhóm đã viết một đoạn mã giúp tạo một mô hình mạng SSD rỗng, sau đó gán tham số từ mô hình mạng VGG16 vừa thu được sang mô hình mạng SSD mới. Ta bỏ qua các lớp liên kết đầy đủ do phần này đã được lược bỏ, lớp pooling cũng được bỏ qua do lớp này không có tham số.

Sau khi đã thu được mô hình mạng SSD, ta có thể bước tiếp đến công đoạn huấn luyện mạng SSD.

3 Xây dựng bản thử nghiệm

Khi đã thu được mô hình hoàn thiện của SSD, nhóm xây dựng một bản thử nghiệm bằng python. Giao diện khá đơn giản, người dùng chỉ việc chọn thư mục chứa ảnh, sau đó chọn ảnh cần nhận diện. Chương trình sẽ xử lý và đưa ra 3 ảnh gồm ảnh gốc ban đầu, ảnh thể hiện các ô bọc mà mạng dự đoán và ảnh biểu thức toán học được tạo ra từ đoạn mã Latex.



Hình 7: Giao diện bản thử nghiệm

4 Đánh giá ưu, nhược điểm

4.1 Ưu điểm

- Hệ thống có thể nhận diện được ký tự ở nhiều kích thước khác nhau.
- Nhận diện được nhiều biểu thức mà nếu chỉ dùng kỹ thuật phân vùng thì khó có thể làm được (ví dụ như biểu thức có ký tự dính nhau).
- Số lượng ký tự nhận diện được khá đa dạng

4.2 Nhược điểm

- Hệ thống có thể nhận diện được ký tự ở nhiều kích thước khác nhau.
- Nhận diện được nhiều biểu thức mà nếu chỉ dùng kỹ thuật phân vùng thì khó có thể làm được (ví dụ như biểu thức có ký tự dính nhau).
- Số lượng ký tự nhận diện được khá đa dạng

Chương 6 Tổng kết

1 Công việc đã làm được

1.1 Chuẩn bị cơ sở lý thuyết

- Hiểu được luận văn của nhóm sinh viên khóa 2011 ở mức độ hiện thực lại được.
- Tìm hiểu về mạng CNN ở mức độ hiểu được quy tắc truyền xuôi và truyền ngược.
- Đã tìm hiểu được cách hiện thực mạng CNN của caffe.
- Đọc và hiểu được các phương pháp giải quyết cho vấn đề nhận dạng biểu thức toán học được nêu trong các công trình nghiên cứu gần đây nhất.

2 Nhận xét bản thử nghiệm hiện tại

Một số vấn đề còn tồn đọng:

- Hệ thống không thể nhận diện các ảnh có nhiễu lớn.
- Hệ thống không thể nhận diện các ảnh có quá nhiều ký tự.

Đây cũng là vấn đề chưa giải quyết được trong luận văn của nhóm sinh viên khóa 2011.



3 Kế hoạch trong giai đoạn tới

Chương 7 Kết luận

1 Kết luận

2 Hướng phát triển trong tương lai

Hệ thống vẫn còn nhiều điểm có thể cải thiện thêm, trong tương lai, có thể:

- Thêm phần hậu xử lý ở các giai đoạn nhận diện ký tự và sinh cây Lexed - BST.
- Tăng thêm số lượng ký tự có thể nhận diện.
- Cải thiện khả năng nhận diện các ký tự có tỉ lệ diện mạo đặc biệt

Tài liệu

- [1] W. He, Y. Luo, and F. Yin, “Context-aware mathematical expression recognition: An end-to-end framework and a benchmark,” pp. 3235–3240, Dec. 2016.
- [2] A. N. Quoc and K. N. Anh, *Nhận dạng biểu thức toán học*. 2015.
- [3] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *CoRR*, vol. abs/1409.1556, 2014.
- [4] Y. Lecun, L. Boutou, and Y. Bengio, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE*, vol. 88, no. 11, pp. 2278 – 2324, Nov. 1998.
- [5] R. H. J. Ha and I. Phillips, “Understanding mathematical expressions from document images,” 1995.
- [6] N. Okamoto and M. Bin, “Recognition of mathematical expressions by using the layout structures of symbols,” 1991.
- [7] R. Zanibbi, D. Blostein, and J. Cordy, “Recognizing mathematical expressions using tree transformation,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 11, pp. 1455 – 1467, Nov. 2002.
- [8] K.-F. Chan and D.-Y. Yeung, “Mathematical expression recognition: a survey,” *International Journal on Document Analysis and Recognition*, vol. 3, no. 1, pp. 3–15, Aug. 2000.
- [9] R. H. Anderson, “Syntax-directed recognition of hand-printed two-dimensional mathematics,” *Symposium on Interactive Systems for Experimental Applied Mathematics: Proceedings of the Association for Computing Machinery Inc. Symposium*, pp. 436–459, Aug. 1967.