

1 Mở đầu

2 Công trình liên quan

3 Phương pháp

4 Kết quả và Tổng kết

5 Tài liệu tham khảo

1 Mở đầu

2 Công trình liên quan

3 Phương pháp

4 Kết quả và Tổng kết

5 Tài liệu tham khảo

Mở đầu I

Mở đầu

Công trình
liên quan

Phương
pháp

Kết quả và
Tổng kết

Tài liệu
tham khảo

Phân vùng ảnh dựa trên ngữ nghĩa là một lĩnh vực quan trọng trong thị giác máy tính và được ứng dụng rộng rãi trong các lĩnh vực như phân tích ảnh y khoa, hệ thống xe tự hành, giám sát video và thực tế tăng cường. Mục tiêu của bài toán này là gán nhãn ngữ nghĩa cho từng điểm ảnh, từ đó phân biệt và xác định các vật thể và đối tượng. Cụ thể, nó sẽ thực hiện việc gán một tập nhãn gồm các loại đối tượng (ví dụ: con người, xe cộ, cây cối, bầu trời, đường xá...) cho tất cả các điểm ảnh trong ảnh. Điều này giúp nội dung hình ảnh được cấp chi tiết hơn so với việc chỉ dự đoán một nhãn cho toàn bộ hình ảnh.

Mở đầu II

Mở đầu

Công trình
liên quan

Phương
pháp

Kết quả và
Tổng kết

Tài liệu
tham khảo

Nhiều phương pháp phân vùng ảnh dựa trên ngữ nghĩa đã được ra đời như phân ngưỡng (thresholding), nhóm dựa trên histogram (histogram-based bundling), lan vùng (region growing), phân cụm K-means (K-means clustering)... Tuy nhiên, nhóm phương pháp truyền thống này thường gặp nhiều hạn chế về độ chính xác khi đối tượng có hình dạng phức tạp hoặc ảnh có độ nhiễu cao (ví dụ: phương pháp phân ngưỡng có thể gặp khó khăn trong việc phân biệt các đối tượng có độ tương phản thấp), khả năng tổng quát hóa kém do thường phụ thuộc nhiều vào đặc điểm thủ công được thiết kế cho từng bài toán cụ thể (ví dụ: phân cụm K-means yêu cầu số lượng cụm được xác định trước) và thời gian tính toán lớn khi phải xử lý ảnh có độ phân giải cao.

Mở đầu III

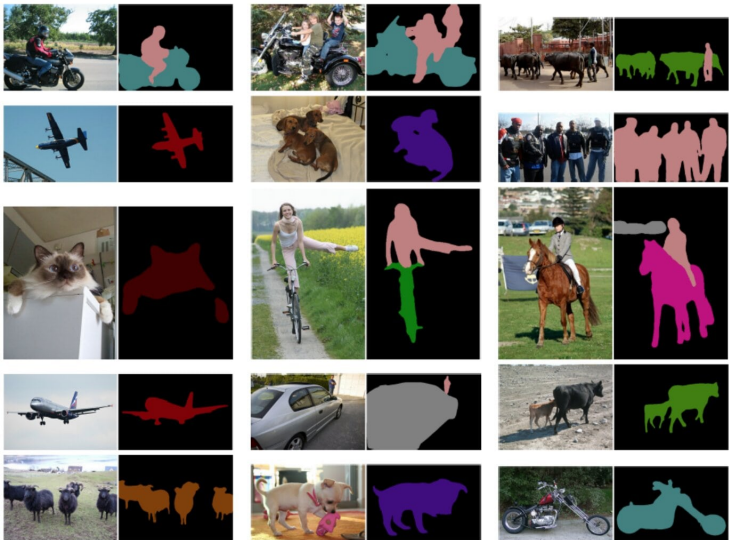
Mở đầu

Công trình liên quan

Phương pháp

Kết quả và Tổng kết

Tài liệu tham khảo



Hình 1: Kết quả phân vùng của DeepLabV3 trên ảnh mẫu [Minaee et al., 2020].

Mở đầu IV

Mở đầu

Công trình
liên quan

Phương
pháp

Kết quả và
Tổng kết

Tài liệu
tham khảo

Tuy nhiên, sự ra đời của các mô hình học sâu, đặc biệt là các mạng nơ-ron tích chập (Convolutional Neural Networks - CNN), đã mang lại những cải tiến hiệu suất đáng chú ý, từ đó dẫn đến sự thay đổi mô hình trong lĩnh vực phân đoạn hình ảnh này. Ví dụ, mô hình DeepLabv3 đã đạt được độ chính xác cao và hiệu suất vượt trội so với các phương pháp truyền thống (Hình 1).

Mở đầu

Công trình
liên quan

Mạng Tích
chập Đầy đủ
(FCN)

Bộ giải
mã-Bộ mã
hoá

U-Net

Kết hợp
thêm cổng
Attention

Phương
pháp

Kết quả và
Tổng kết

Tài liệu
tham khảo

1 Mở đầu

2 Công trình liên quan

- U-Net
- Kết hợp thêm cổng Attention

3 Phương pháp

4 Kết quả và Tổng kết

5 Tài liệu tham khảo

Mạng Tích chập Dây đủ (FCN) I

Mở đầu

Công trình liên quan

Mạng Tích chập Dây đủ (FCN)

Bộ giải mã-Bộ mã hoá

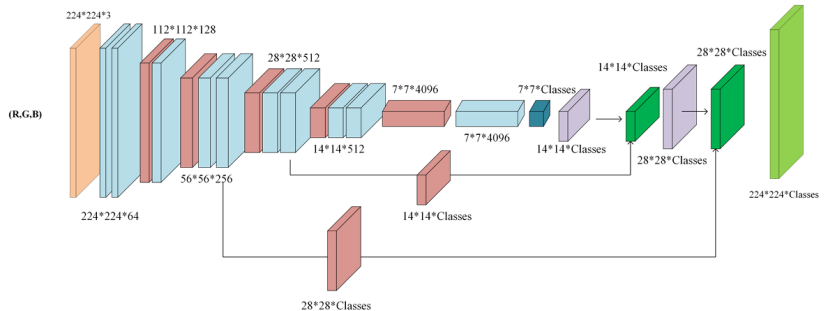
U-Net

Kết hợp thêm công Attention

Phương pháp

Kết quả và Tổng kết

Tài liệu tham khảo



Hình 2: Kiến trúc của FCN-8 [Piramanayagam et al., 2018].

Mạng Tích chập Dây đủ (FCN) II

Mở đầu

Công trình
liên quan

Mạng Tích
chập Dây đủ
(FCN)

Bộ giải
mã-Bộ mã
hoá

U-Net

Kết hợp
thêm công
Attention

Phương
pháp

Kết quả và
Tổng kết

Tài liệu
tham khảo

Mạng tích chập toàn bộ (Fully Convolution Network, gọi tắt là FCN) [Long et al., 2015] là mạng toàn bộ các lớp đều là lớp Tích chập (CNN).

FCN đã đưa ra một ý tưởng quan trọng là “Kết nối Ngắt quãng” (Skip Connection), trong đó ánh xạ đặc trưng (feature map) của một lớp sẽ được mở rộng mẫu (up-sampled) và sau đó kết hợp với ánh xạ đặc trưng của một lớp khác (ở hình 2).

Việc dùng kết nối ngắt quãng này sẽ kết hợp được thông tin giữa các lớp với nhau và đưa ra cách phân vùng chính xác. FCN được xem là một cột mốc quan trọng trong tác vụ phân vùng ảnh dựa trên ngữ nghĩa thể nhưng FCN vẫn còn một số hạn chế như không đủ nhanh để có thể suy diễn trong thời gian thực và không xử lý ngữ cảnh toàn cục của ảnh một cách hiệu quả.

Bộ giải mã-Bộ mã hoá I

Mở đầu

Công trình liên quan

Mạng Tích chập Dây chuyền (FCN)

Bộ giải mã-Bộ mã hoá

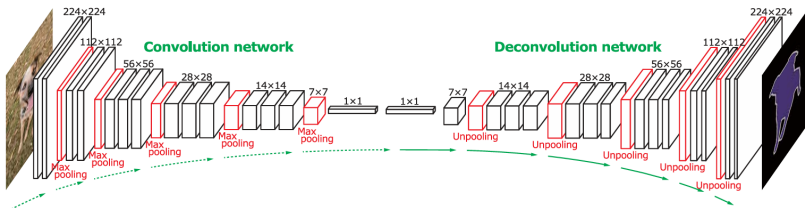
U-Net

Kết hợp thêm công Attention

Phương pháp

Kết quả và Tổng kết

Tài liệu tham khảo



Hình 3: Kiến trúc mạng được đề xuất bởi Noh và cộng sự [Noh et al., 2015].

Bộ giải mã-Bộ mã hoá II

Mở đầu

Công trình
liên quan

Mạng Tích
chập Dây đủ
(FCN)

Bộ giải
mã-Bộ mã
hoá

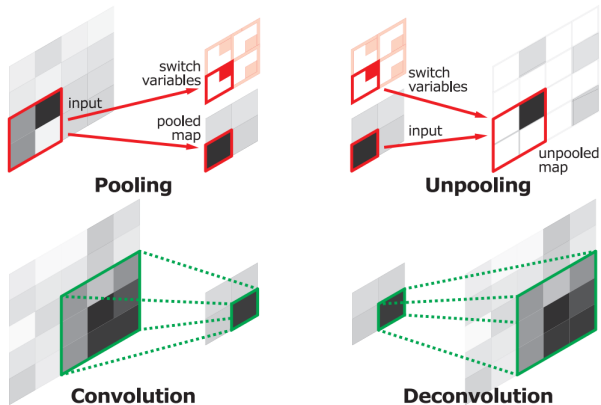
U-Net

Kết hợp
thêm công
Attention

Phương
pháp

Kết quả và
Tổng kết

Tài liệu
tham khảo



Hình 4: Ví dụ về Tích chập chuyển vị (Deconvolution) và Mở gộp (Unpool) trong kiến trúc Bộ mã hoá-Bộ giải mã (hình được lấy từ [Zhang et al., 2017]).

Bộ giải mã-Bộ mã hoá III

Mở đầu

Công trình
liên quan

Mạng Tích
chập Đầy đủ
(FCN)

Bộ giải
mã-Bộ mã
hoá

U-Net

Kết hợp
thêm công
Attention

Phương
pháp

Kết quả và
Tổng kết

Tài liệu
tham khảo

Noh và cộng sự [Noh et al., 2015] đã đưa ra một trong những mô hình đầu tiên sử dụng cơ chế gọi là “Bộ mã hoá-Bộ giải mã” (Encoder-Decoder). Trong đó mạng sẽ được chia ra làm hai phần, phần được gọi là “Bộ mã hoá” và phần được gọi là “Bộ giải mã”.

Trong nghiên cứu của Noh và cộng sự [Noh et al., 2015], Bộ mã hoá là một mạng con gồm nhiều lớp tích chập và lớp gộp (pool), còn Bộ giải mã là một mạng con gồm nhiều lớp để giải mã lớp tích chập và lớp gộp ở bộ mã hoá, gọi là Lớp tích chập chuyển vị (transposed convolution hay deconvolution) (ở hình 4) và Lớp mở gộp (unpool).

U-Net I

Mở đầu

Công trình liên quan

Mạng Tích chập Dây đủ (FCN)

Bộ giải mã-Bộ mã hoá

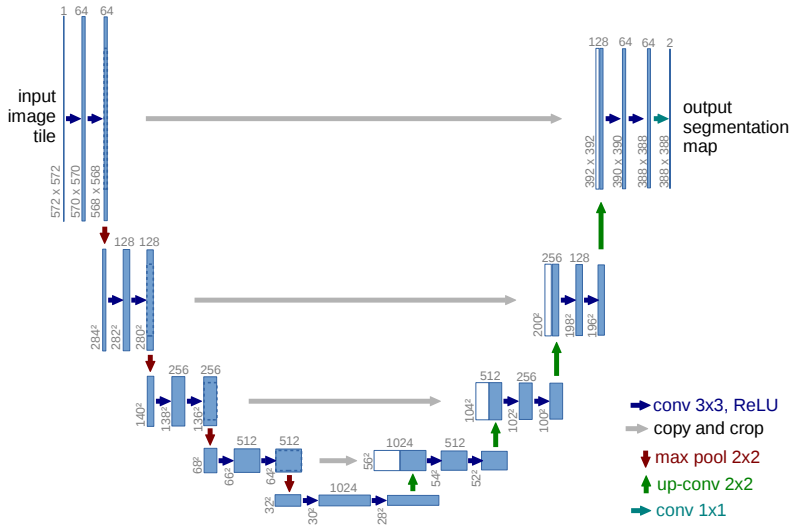
U-Net

Kết hợp thêm công Attention

Phương pháp

Kết quả và Tổng kết

Tài liệu tham khảo



Hình 5: Mô hình U-Net gốc của Ronneberger và cộng sự [Ronneberger et al., 2015]

Mở đầu

Công trình
liên quan

Mạng Tích
chập Đầy đủ
(FCN)

Bộ giải
mã-Bộ mã
hoá

U-Net

Kết hợp
thêm công
Attention

Phương
pháp

Kết quả và
Tổng kết

Tài liệu
tham khảo

Nhờ nghiên cứu của Noh [Noh et al., 2015] và Long [Long et al., 2015], Ronneberger và cộng sự [Ronneberger et al., 2015] đã đưa ra mô hình U-Net, mô hình này kết hợp hai ý tưởng chính trên, “Kết nối ngắt quãng” và “Bộ mã hoá-Bộ giải mã”. Ở kết nối ngắt quãng (mũi tên màu xám ở hình 5), mô hình U-Net nối hai ánh xạ đặc trưng lại với nhau và sau đó cho một lớp tích chập cùng với một hàm kích hoạt không tuyến tính.

Mô hình U-Net chia làm hai phần. Phần lấy mẫu xuống (down-sampling) hay được xem như bộ mã hoá (phần bên trái của hình 5) có nhiệm vụ trích xuất các đặc trưng của ảnh, phần này được cấu tạo từ các lớp tích chập có kích thước 3×3 . Phần còn lại là phần mở rộng mẫu (up-sampling) hay bộ giải mã (phần bên phải của hình 5) sẽ làm giảm số lượng kênh của các ánh xạ đặc trưng, đồng thời các ánh xạ đặc trưng ở phần mở rộng được nối với ánh xạ bên phần lấy mẫu xuống (kết nối ngắt quãng), bộ giải mã được tạo thành từ nhiều tích chập chuyển vị như của [Noh et al., 2015] hoặc có thể dùng các lớp nội suy song tuyến (bilinear interpolation).

Kết hợp thêm cổng Attention I

Mở đầu

Công trình
liên quan

Mạng Tích
chập Đầy đủ
(FCN)

Bộ giải
mã-Bộ mã
hoá

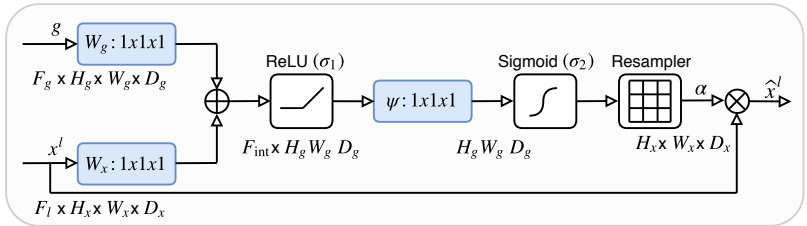
U-Net

Kết hợp
thêm cổng
Attention

Phương
pháp

Kết quả và
Tổng kết

Tài liệu
tham khảo



Hình 6: Cổng attention trong bài báo của Oktay và cộng sự [Oktay et al., 2018].

Kết hợp thêm cổng Attention II

Mở đầu

Công trình liên quan

Mạng Tích chấp Dầy đủ (FCN)

Bộ giải mã-Bộ mã hoá

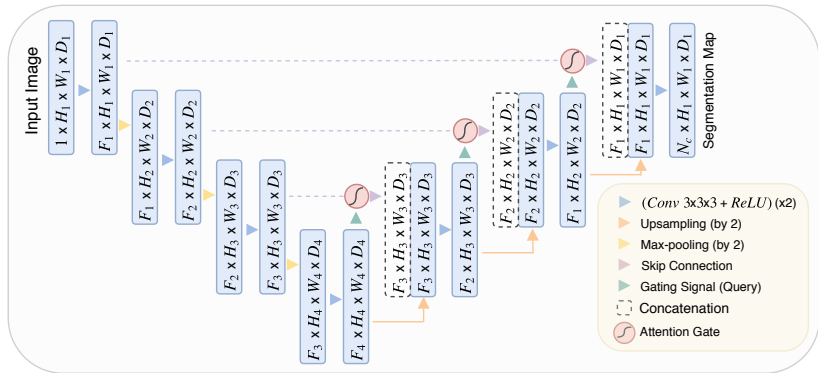
U-Net

Kết hợp thêm cổng Attention

Phương pháp

Kết quả và Tổng kết

Tài liệu tham khảo



Hình 7: Mô hình Attention U-Net trong bài báo của Oktay và cộng sự [Oktay et al., 2018].

Kết hợp thêm cổng Attention III

Mở đầu

Công trình
liên quan

Mạng Tích
chập Đầy đủ
(FCN)

Bộ giải
mã-Bộ mã
hoá

U-Net

Kết hợp
thêm cổng
Attention

Phương
pháp

Kết quả và
Tổng kết

Tài liệu
tham khảo

Để cải tiến U-Net, nhóm tác giả đưa ra cải tiến mà nhóm tác giả thấy phù hợp và dễ hiểu nhất. Đó là thêm cơ chế Cổng Attention (Attention Gate) vào đoạn “Kết nối ngắn quãng” của U-Net gốc. Cải tiến này được dựa trên bài báo của Oktay và cộng sự [Oktay et al., 2018]. Thay vì chỉ nối hai ánh xạ đặc trưng lại với nhau như trong bài báo [Ronneberger et al., 2015], cải tiến này sẽ kết hợp hai ánh xạ đặc trưng ấy thông qua một cổng Attention (như hình 6) sau đó mới tiến hành nối lại. Toàn bộ mô hình U-Net cộng với cải tiến trong bài báo gốc [Oktay et al., 2018] ở hình 7.

Mục lục

Mở đầu

Công trình
liên quan

Phương
pháp

Tập dữ liệu

Mô hình
U-Net

Mô hình
U-Net với cơ
chế
Attention

Huấn luyện

Kết quả và
Tổng kết

Tài liệu
tham khảo

1 Mở đầu

2 Công trình liên quan

3 Phương pháp

- Tập dữ liệu
- Mô hình U-Net
- Mô hình U-Net với cơ chế Attention
- Huấn luyện

4 Kết quả và Tổng kết

5 Tài liệu tham khảo

Tập dữ liệu I

Mở đầu

Công trình
liên quan

Phương
pháp

Tập dữ liệu

Mô hình
U-Net

Mô hình
U-Net với cơ
chế
Attention

Huấn luyện

Kết quả và
Tổng kết

Tài liệu
tham khảo



(a) Ảnh ví dụ



(b) Nhãn tương ứng

Hình 8: Một ví dụ về ảnh và nhãn tương ứng của nó (đã được gộp lại còn 8 nhãn) trong bộ dữ liệu Cityscapes [Cordts et al., 2016].

Tập dữ liệu II

Mở đầu

Công trình
liên quan

Phương
pháp

Tập dữ liệu

Mô hình
U-Net

Mô hình
U-Net với cơ
chế
Attention

Huấn luyện

Kết quả và
Tổng kết

Tài liệu
tham khảo

Nhóm tác giả sử dụng bộ dữ liệu Cityscapes [Cordts et al., 2016] cho việc huấn luyện, thử nghiệm và đánh giá.

Bộ dữ liệu ban đầu sẽ bao gồm 2975 ảnh cho tập huấn luyện, 500 ảnh cho tập thẩm định với nhãn được công khai và 1525 ảnh thử nghiệm với nhãn được giấu đi.

Tập dữ liệu ban đầu sẽ gồm 34 nhãn được dùng cho tác vụ phân vùng ảnh dựa trên đối tượng, để làm phù hợp cho tác vụ phân vùng ảnh dựa trên ngữ nghĩa, nhóm tác giả đã gộp lại còn 8 nhãn tất cả. Một ví dụ về bộ dữ liệu này được đưa ra ở hình 8.

Mô hình U-Net I

Mở đầu

Công trình liên quan

Phương pháp

Tập dữ liệu

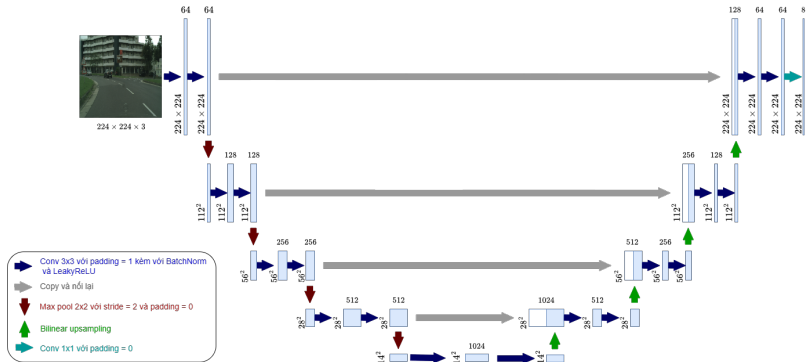
Mô hình U-Net

Mô hình U-Net với cơ chế Attention

Huấn luyện

Kết quả và Tổng kết

Tài liệu tham khảo



Hình 9: Mô hình U-Net của nhóm tác giả (đã được thay đổi một chút so với bản gốc [Ronneberger et al., 2015] và được sử dụng trên ảnh màu RGB có kích thước 224×224). Mỗi hình chữ nhật màu xanh tương ứng với một ánh xạ đặc trưng nhiều kênh. Số kênh của ánh xạ ấy được kí hiệu phía trên hình chữ nhật. Hình chữ nhật màu trắng tương ứng với một bản sao của hình chữ nhật màu xanh phía sau mũi tên.

Mô hình U-Net II

Mở đầu

Công trình
liên quan

Phương
pháp

Tập dữ liệu

Mô hình
U-Net

Mô hình
U-Net với cơ
chế
Attention

Huấn luyện

Kết quả và
Tổng kết

Tài liệu
tham khảo

Mô hình đầu tiên mà nhóm tác giả chọn để thực hiện tác vụ này là mô hình U-Net [Ronneberger et al., 2015]. Mô hình được trực quan hoá trong hình 5.

Vẫn giữ nguyên so với bài báo gốc, nhóm tác giả chọn số lượng bộ lọc của lớp tích chập ban đầu là 64. Đầu tiên ở lớp Tích chập với nhân kích thước 3×3 (Conv 3x3 trong hình 5), nhóm tác giả chọn tham số padding là 1 thay vì 0 như trong bài báo gốc, điều này giúp tránh trường hợp ảnh có kích thước lẻ khi đi qua lớp Gộp Tối đa với nhân kích thước 2×2 (Max pool 2x2 trong hình 9) cũng như là làm cho ánh xạ đặc trưng đầu ra có cùng kích thước với ảnh đầu vào.

Mô hình U-Net III

Mở đầu

Công trình
liên quan

Phương
pháp

Tập dữ liệu

Mô hình
U-Net

Mô hình
U-Net với cơ
chế
Attention

Huấn luyện

Kết quả và
Tổng kết

Tài liệu
tham khảo

Tiếp theo, nhóm tác giả chọn hàm kích hoạt là Hàm chỉnh lưu rò rỉ (LeakyReLU) [Xu et al., 2015] thay vì Hàm chỉnh lưu thông thường (ReLU) như trong bài báo gốc, bởi vì sau nhiều lần thử nghiệm nhóm tác giả thấy Hàm chỉnh lưu rò rỉ phù hợp hơn với dữ liệu và cho kết quả tốt hơn.

Để cho quá trình huấn luyện nhanh và ổn định hơn, ở mỗi lớp Tích chập, trước khi qua hàm kích hoạt ReLU thì sẽ đi qua một lớp Chuẩn hoá Hàng loạt (BatchNorm trong hình 9) [Ioffe and Szegedy, 2015]. Cuối cùng, ở đoạn lấy mẫu lên, nhóm tác giả chọn cách dùng nội suy song tuyến (như trong hình 9).

Mô hình U-Net với cơ chế Attention I

Mở đầu

Công trình liên quan

Phương pháp

Tập dữ liệu

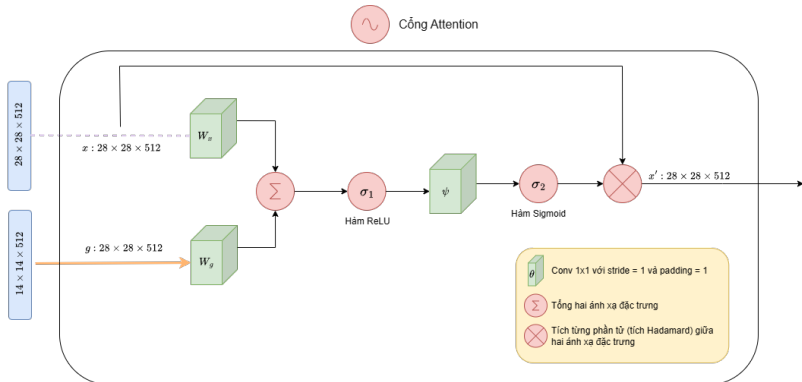
Mô hình U-Net

Mô hình U-Net với cơ chế Attention

Huấn luyện

Kết quả và Tổng kết

Tài liệu tham khảo



Hình 10: Cổng Attention của nhóm tác giả (đã thay đổi một chút so với bản gốc [Oktay et al., 2018]).

Mô hình U-Net với cơ chế Attention II

Mở đầu

Công trình
liên quan

Phương
pháp

Tập dữ liệu

Mô hình
U-Net

Mô hình
U-Net với cơ
chế
Attention

Huấn luyện

Kết quả và
Tổng kết

Tài liệu
tham khảo

Để làm đơn giản đi cổng Attention, nhóm tác giả đã bỏ đi Bộ lấy mẫu lại (Resampler) trong bài báo gốc [Oktay et al., 2018] (hoặc trong hình 6). Mô hình U-Net kèm cơ chế Attention được nhóm tác giả đơn giản hoá và trực quan hoá ở hình 11.

Ở cổng Attention, hình hộp màu xanh là một lớp tích chập có kích thước 1×1 với bộ tham số của nó (ví dụ như bộ tham số W_g ở hình 10). Khi ánh xạ đặc trưng ở lớp trước đó, gọi là x và ánh xạ đặc trưng ở lớp hiện tại sau khi được lấy mẫu lên, gọi là g cùng đi qua cổng Attention, đầu tiên chúng sẽ đi qua lớp tích chập của tương ứng với mình W_x cho x và W_g cho g , riêng lớp tích chập của g sẽ có bias cộng thêm, gọi là b_g . Sau khi cả hai đã đi qua lớp tích chập của mình, chúng sẽ được cộng lại với nhau và đi qua một hàm kích hoạt σ_1 (ở đây là hàm ReLU), gọi kết quả của đoạn này là q .

Mô hình U-Net với cơ chế Attention III

Mở đầu

Công trình
liên quan

Phương
pháp

Tập dữ liệu

Mô hình
U-Net

Mô hình
U-Net với cơ
chế
Attention

Huấn luyện

Kết quả và
Tổng kết

Tài liệu
tham khảo

$$q = \sigma_1 (W_x^T x + W_g^T g + b_g)$$

Hàm kích hoạt σ_1 (ReLU)

Tham số của lớp tích chập mà x đi qua

Tham số của lớp tích chập mà g đi qua

Mô hình U-Net với cơ chế Attention IV

Mở đầu

Công trình
liên quan

Phương
pháp

Tập dữ liệu

Mô hình
U-Net

Mô hình
U-Net với cơ
chế
Attention

Huấn luyện

Kết quả và
Tổng kết

Tài liệu
tham khảo

Sau khi có được q , ta cho q đi qua một lớp tích chập nữa có kích thước 1×1 và có bộ tham số là ψ , tương tự W_g , lớp này cũng có bias và ta gọi bias đó là b_ψ . Tiếp theo đó cho qua hàm kích hoạt σ_2 (ở đây là hàm Sigmoid), sau đó ta lấy giá trị đã kích hoạt này nhân với từng phần tử của x ban đầu. Cuối cùng, ta được giá trị đầu ra của cổng Attention, gọi là x' .

$$x' = \sigma_2(\psi^T q + b_\psi) \odot x$$

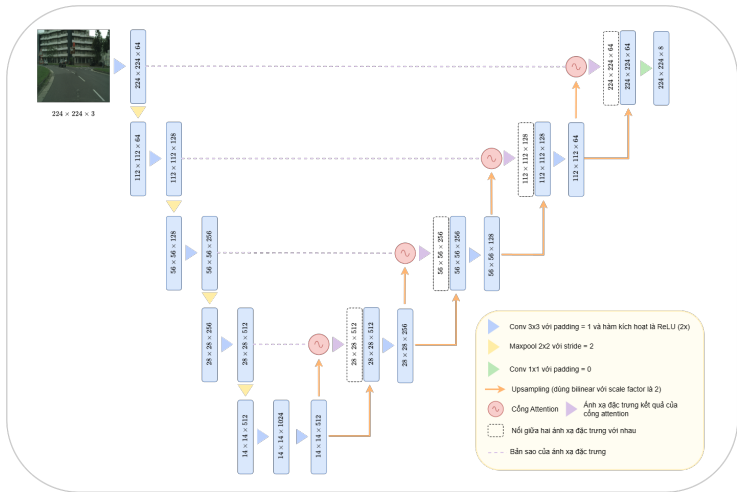
Hàm kích hoạt σ_2 (Sigmoid)

Tích Hadamard hay tích của từng phần tử với nhau

Khi có được giá trị đầu ra x' , ta sẽ tiến hành nối lại với ánh xạ đặc trưng của lớp hiện tại tương tự như trong bài báo U-Net gốc [Ronneberger et al., 2015] (trực quan hoá ở hình 11).

Mô hình U-Net với cơ chế Attention V

Mở đầu
Công trình liên quan
Phương pháp
Tập dữ liệu
Mô hình U-Net
Mô hình U-Net với cơ chế Attention
Huấn luyện
Kết quả và Tổng kết
Tài liệu tham khảo



Hình 11: Mô hình U-Net kèm cơ chế Attention của nhóm tác giả (đã được thay một chút so với bản gốc [Okta et al., 2018] và được sử dụng trên ảnh màu RGB có kích thước 224 × 224).

Huấn luyện I

Mở đầu

Công trình
liên quan

Phương
pháp

Tập dữ liệu

Mô hình
U-Net

Mô hình
U-Net với cơ
chế
Attention

Huấn luyện

Kết quả và
Tổng kết

Tài liệu
tham khảo

Cả ba mô hình đều sử dụng chung một hàm mất mát là hàm Entropy Chéo (Cross Entropy). Sau đó hàm mất mát này sẽ được tối ưu bằng thuật toán AdamW [Loshchilov and Hutter, 2017] trong đó tỉ lệ học (learning rate) được thiết lập là 0.01 và độ phân rã trọng số (weight decay) được thiết lập là 1×10^{-5} .

```
optimizer = torch.optim.Adam(model.parameters(),  
                               lr=0.01,  
                               weight_decay = 1e-5)  
loss = nn.CrossEntropyLoss()
```

Ngoài ra mỗi mô hình được huấn luyện trên 25 epoch (do giới hạn về phần cứng) và có kích thước batch là 32. Để đánh giá được độ tốt của mô hình, nhóm tác giả sử dụng độ đo MeanIoU, với IoU được chọn là độ đo Jaccard (hay chỉ số Jaccard) cho nhiều lớp [Li et al., 2015].

Huấn luyện II

Mở đầu

Công trình
liên quan

Phương
pháp

Tập dữ liệu

Mô hình
U-Net

Mô hình
U-Net với cơ
chế
Attention

Huấn luyện

Kết quả và
Tổng kết

Tài liệu
tham khảo

Ngoài ra, việc khởi tạo giá trị ban đầu cho các trọng số cũng cực kì quan trọng, trong nghiên cứu này, nhóm tác giả chọn khởi tạo các giá trị ấy từ một phân phối chuẩn có trung bình là 0 và phương sai là $2/(N \times M)$ với N là kích thước hiện tại của lớp Tích chập và M là số kênh đầu ra của lớp Tích chập ấy. Ví dụ, một lớp Tích chập có kích thước 3×3 và số kênh đầu ra là 64 thì $N = 9$ và $M = 64$, riêng giá trị bias sẽ được khởi tạo giá trị ban đầu là 0.

```
def init_weight(m):  
    if isinstance(m, nn.Conv2d):  
        N = m.kernel_size[0]*m.kernel_size[1]*m.out_channels  
        torch.nn.init.normal_(m.weight, 0.0, 2/N)  
        if m.bias is not None:  
            torch.nn.init.zeros_(m.bias)
```

Mục lục

Mở đầu

Công trình
liên quan

Phương
pháp

Kết quả và
Tổng kết

Kết quả
Tổng kết

Tài liệu
tham khảo

1 Mở đầu

2 Công trình liên quan

3 Phương pháp

4 Kết quả và Tổng kết

- Kết quả
- Tổng kết

5 Tài liệu tham khảo

Mở đầu

Công trình
liên quan

Phương
pháp

Kết quả và
Tổng kết

Kết quả
Tổng kết

Tài liệu
tham khảo

Bảng 1: So sánh giữa hai mô hình U-Net và Attention U-Net của nhóm tác giả bằng MeanIoU cao nhất và thời gian huấn luyện

Mô hình	MeanIoU cao nhất	Thời gian huấn luyện
U-Net (hình 9)	0.714	2 giờ
Attention U-Net (hình 11)	0.74	2.5 giờ

Mở đầu

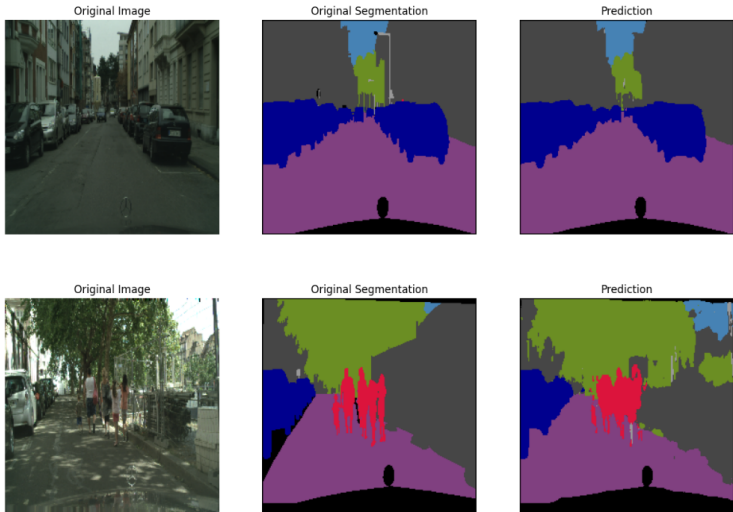
Công trình
liên quan

Phương
pháp

Kết quả và
Tổng kết

Kết quả
Tổng kết

Tài liệu
tham khảo



Hình 12: Dự đoán của Attention U-Net

Kết quả III

Mở đầu

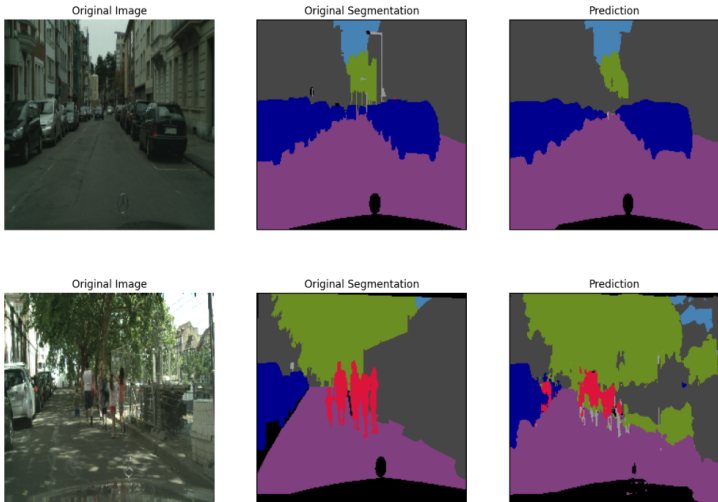
Công trình
liên quan

Phương
pháp

Kết quả và
Tổng kết

Kết quả
Tổng kết

Tài liệu
tham khảo



Hình 13: Dự đoán của mô hình U-Net

Kết quả IV

Mở đầu

Công trình
liên quan

Phương
pháp

Kết quả và
Tổng kết

Kết quả
Tổng kết

Tài liệu
tham khảo

Có thể thấy, việc sử dụng thêm cổng Attention cho phần Kết nối ngắt quãng đã cho kết quả tốt hơn rất nhiều khi MeanIoU cao nhất tăng thêm gần 0.03, thế nhưng thời gian huấn luyện sẽ tăng lên thêm 0.5 giờ (theo bảng 1), ngoài ra theo hình 13 và 12 thì dự đoán của Attention U-Net phần nào đó chuẩn xác và gần với nhãn gốc hơn so với U-Net.

Tổng kết

Mở đầu

Công trình
liên quan

Phương
pháp

Kết quả và
Tổng kết

Kết quả
Tổng kết

Tài liệu
tham khảo

Có thể thấy nhóm tác giả đã áp dụng thành công hai mô hình U-Net [Ronneberger et al., 2015] và cải tiến của nó là Attention U-Net [Oktay et al., 2018] vào bài toán phân vùng ảnh bằng ngữ nghĩa trên tập dữ liệu Cityscapes [Cordts et al., 2016], kết quả cho ra ngoài sức mong đợi của nhóm tác giả, việc áp dụng tuân thủ các quy tắc của bài báo gốc và thêm sự thay đổi của nhóm tác giả đã góp phần không nhỏ cho kết quả này. Thế nhưng, trong tương lai, nhóm tác giả sẽ sử dụng các mô hình khác mạnh mẽ và hiệu quả hơn ví dụ như U-Net++ [Zhou et al., 2018], đây là mô hình mà nhóm tác giả đã xem xét sử dụng, nhưng vì kiến trúc phức tạp của nó nên nhóm tác giả đã không chọn, ngoài các mô hình dựa vào U-Net, các mô hình khác mới và hiện đại hơn cũng có thể được xem xét như SegFormer [Xie et al., 2021] hay VLTSeg [Hümmer et al., 2023].

1 Mở đầu

2 Công trình liên quan

3 Phương pháp

4 Kết quả và Tổng kết

5 Tài liệu tham khảo

Tài liệu tham khảo I

Mở đầu

Công trình
liên quan

Phương
pháp

Kết quả và
Tổng kết

Tài liệu
tham khảo



Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., and Schiele, B. (2016).

The cityscapes dataset for semantic urban scene understanding.

CoRR, abs/1604.01685.



He, K., Zhang, X., Ren, S., and Sun, J. (2015).

Deep residual learning for image recognition.

CoRR, abs/1512.03385.



Hümmer, C., Schwonberg, M., Zhou, L., Cao, H., Knoll, A., and Gottschalk, H. (2023).

Vltseg: Simple transfer of clip-based vision-language representations for domain generalized semantic segmentation.



Ioffe, S. and Szegedy, C. (2015).

Batch normalization: Accelerating deep network training by reducing internal covariate shift.

CoRR, abs/1502.03167.



Li, L., Wu, Y., and Ye, M. (2015).

Experimental comparisons of multi-class classifiers.

Informatica, 39(1).

Tài liệu tham khảo II

Mở đầu

Công trình
liên quan

Phương
pháp

Kết quả và
Tổng kết

Tài liệu
tham khảo



Long, J., Shelhamer, E., and Darrell, T. (2015).

Fully convolutional networks for semantic segmentation.

In 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 3431–3440.



Loshchilov, I. and Hutter, F. (2017).

Fixing weight decay regularization in adam.

CoRR, abs/1711.05101.



Minaee, S., Boykov, Y., Porikli, F., Plaza, A., Kehtarnavaz, N., and Terzopoulos, D. (2020).

Image segmentation using deep learning: A survey.

CoRR, abs/2001.05566.



Noh, H., Hong, S., and Han, B. (2015).

Learning deconvolution network for semantic segmentation.

CoRR, abs/1505.04366.



Oktay, O., Schlemper, J., Folgoc, L. L., Lee, M. C. H., Heinrich, M. P., Misawa, K., Mori, K., McDonagh, S. G., Hammerla, N. Y., Kainz, B., Glocker, B., and Rueckert, D. (2018).

Attention u-net: Learning where to look for the pancreas.

CoRR, abs/1804.03999.

Tài liệu tham khảo III

Mở đầu

Công trình
liên quan

Phương
pháp

Kết quả và
Tổng kết

Tài liệu
tham khảo



Piramanayagam, S., Saber, E., Schwartzkopf, W., and Koehler, F. (2018).

Supervised classification of multisensor remotely sensed images using a deep learning framework.

Remote Sensing, 10:1429.



Ronneberger, O., Fischer, P., and Brox, T. (2015).

U-net: Convolutional networks for biomedical image segmentation.

CoRR, abs/1505.04597.



Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J. M., and Luo, P. (2021).

Segformer: Simple and efficient design for semantic segmentation with transformers.

In Ranzato, M., Beygelzimer, A., Dauphin, Y., Liang, P., and Vaughan, J. W., editors, *Advances in Neural Information Processing Systems*, volume 34, pages 12077–12090. Curran Associates, Inc.



Xu, B., Wang, N., Chen, T., and Li, M. (2015).

Empirical evaluation of rectified activations in convolutional network.

CoRR, abs/1505.00853.



Zhang, Z., Liu, Q., and Wang, Y. (2017).

Road extraction by deep residual u-net.

CoRR, abs/1711.10684.

Tài liệu tham khảo IV

Mở đầu

Công trình
liên quan

Phương
pháp

Kết quả và
Tổng kết

Tài liệu
tham khảo



Zhou, Z., Siddiquee, M. M. R., Tajbakhsh, N., and Liang, J. (2018).

Unet++: A nested u-net architecture for medical image segmentation.

CoRR, abs/1807.10165.