



1

Nội dung buổi học

- Giới thiệu về bài toán phân đoạn ảnh
- Phân đoạn ảnh dựa trên ngưỡng
- Phân đoạn ảnh dựa vào cạnh
- Phân đoạn ảnh dựa vào miền
- Phân đoạn ảnh phân cụm
- Kỹ thuật học sâu

2

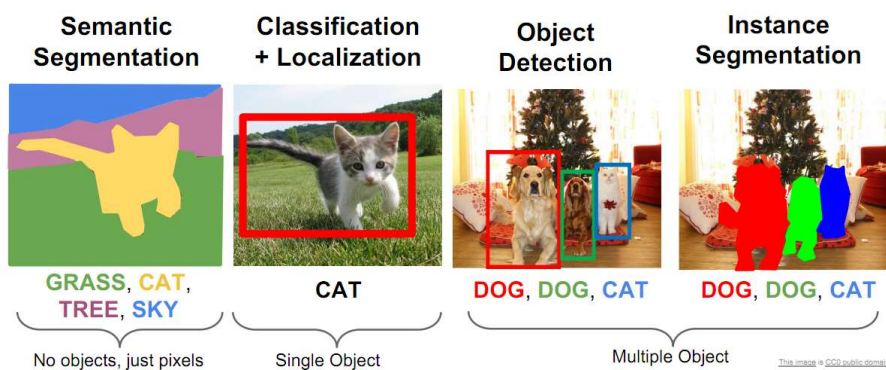
Giới thiệu bài toán phân vùng ảnh



3

3

Các bài toán thị giác máy

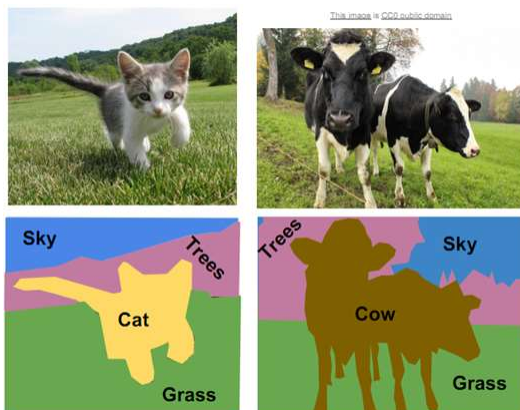


4

4

Phân đoạn ảnh

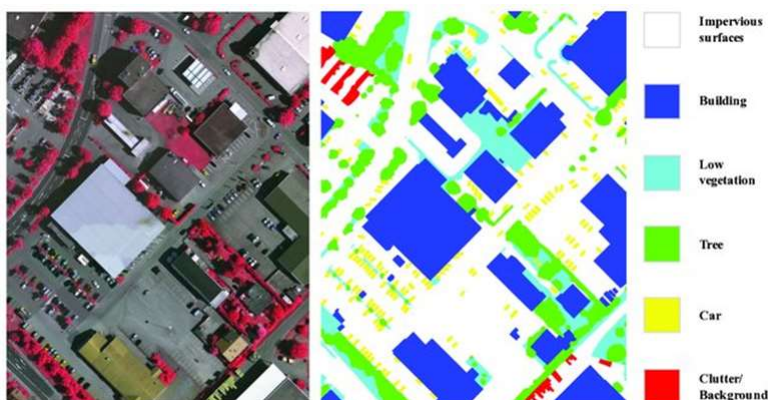
- Phân lớp từng điểm ảnh trong ảnh
- Không phân biệt các đối tượng cùng lớp trong ảnh



5

Một số ứng dụng phân đoạn ảnh

- Phân đoạn ảnh vệ tinh và hàng không



6

Một số ứng dụng phân đoạn ảnh

- Xe tự hành

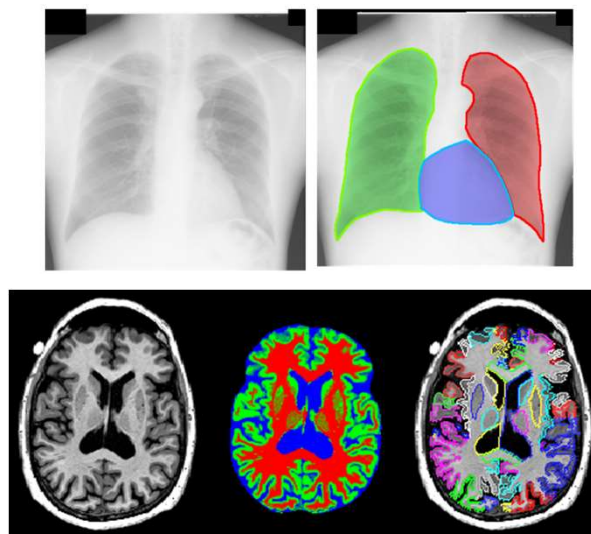


7

7

Một số ứng dụng phân đoạn ảnh

- Y tế



8

8

Một số ứng dụng phân đoạn ảnh

- OCR



figure, table, section heading, caption, list and paragraph



VINBIODATA



VINGROUP

9

9

Phân loại các giải pháp phân đoạn ảnh

- Phân đoạn dựa vào ngưỡng (thresholding)
- Phân đoạn dựa trên phát hiện biên/ cạnh (edge-based)
- Phân đoạn dựa trên vùng (region-based)
- Phân đoạn dựa trên phân cụm (clustering)
- Phân đoạn dựa trên kỹ thuật học sâu



VINBIODATA



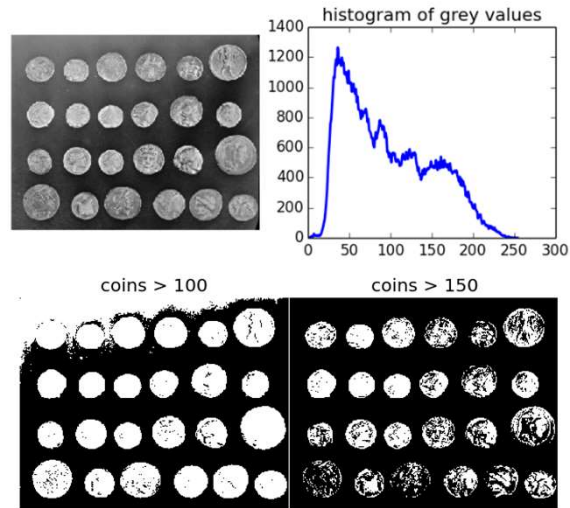
VINGROUP

10

10

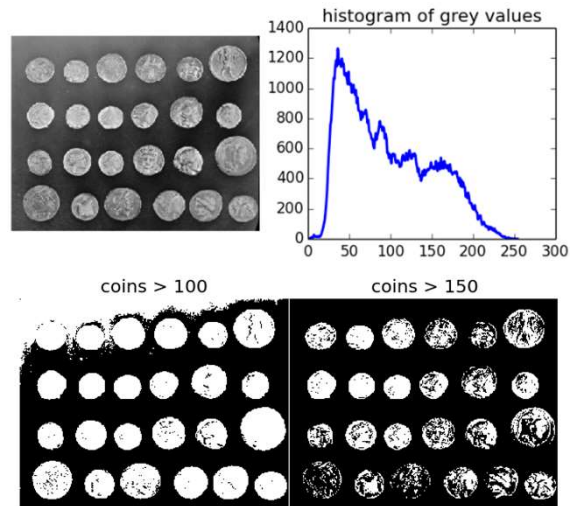
Giải thuật phân đoạn dựa trên ngưỡng

- Phân đoạn ảnh dựa trên ngưỡng đơn → ảnh nhị phân



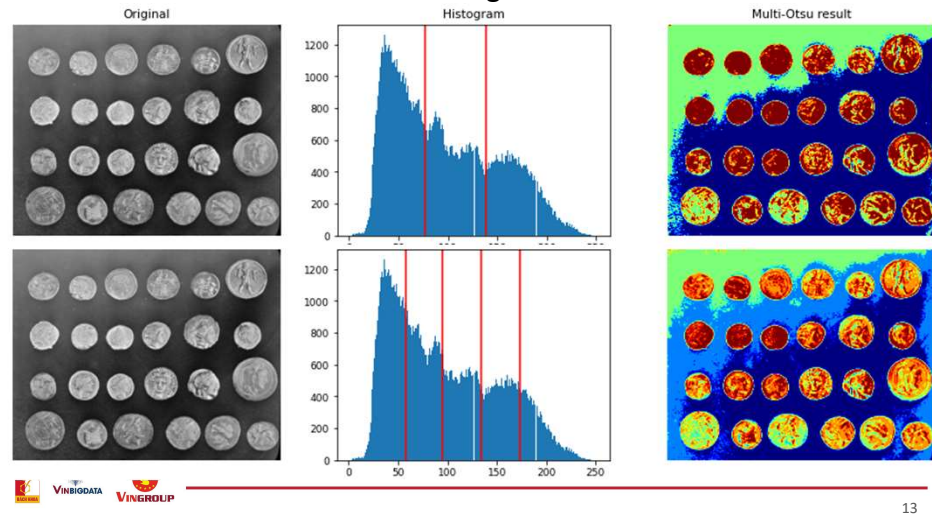
Giải thuật phân đoạn dựa trên ngưỡng

- Phân đoạn ảnh dựa trên ngưỡng đơn → ảnh nhị phân



Giải thuật phân đoạn dựa trên ngưỡng (2)

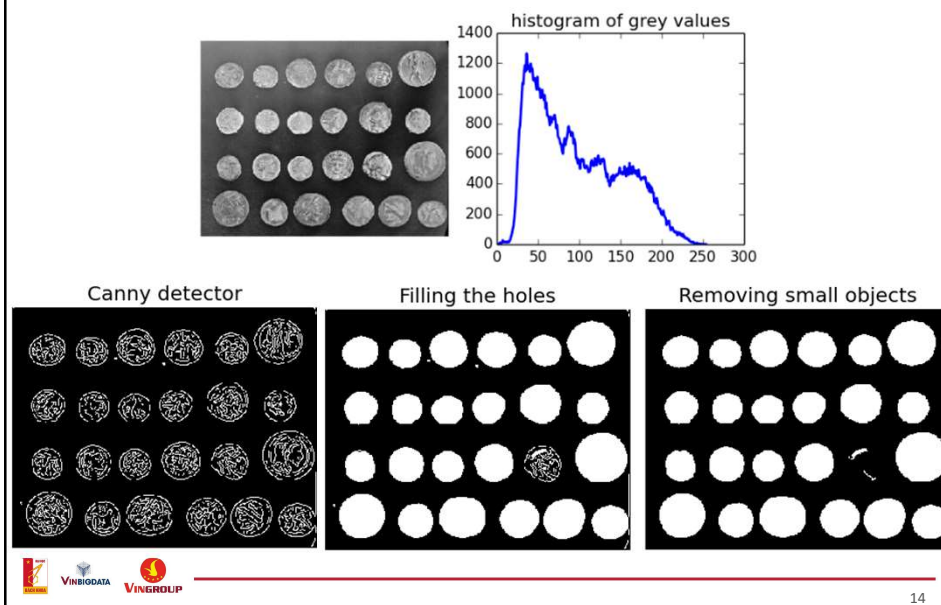
- Phân đoạn ảnh dựa trên đa ngưỡng
 - Vd: Multi-otsu thresholding



13

13

Giải thuật phân đoạn dựa vào biên



14

14

Giải thuật phân đoạn dựa trên vùng

- Lan tỏa/ Phát triển vùng (Region Growing)
- Tách và hợp (Splitting and Merging)

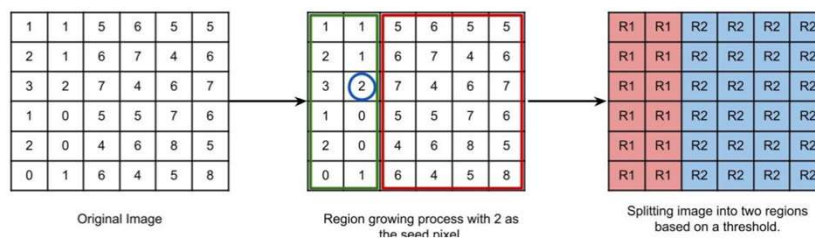


15

15

Lan tỏa vùng (Region Growing)

- B1: Lựa chọn một điểm seed (theo tiêu chí người dùng)
- B2: kiểm tra các điểm lân cận để tìm ra các điểm tương đồng (sai lệch độ sáng nhỏ hơn ngưỡng T), thêm điểm tương đồng vào vùng
- B3: Lặp lại bước 2 với tất cả các điểm trong vùng



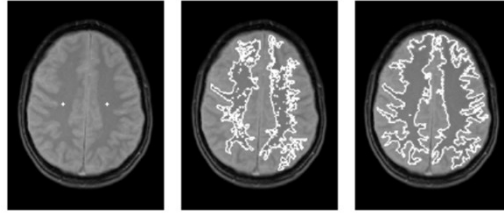
16

16

Lan tỏa vùng (Region Growing)

- Ưu điểm

- Đơn giản, chỉ cần các điểm seed point
- Dễ dàng tách những vùng đồng nhất
- Kết quả tốt với ảnh có biên rõ ràng



Ví dụ lan tỏa vùng từ 2 seed points

- Nhược điểm

- Chi phí tính toán lớn
- Tính toán cục bộ, không phải toàn cục
- Nhạy với nhiễu



17

17

Tách và hợp (Splitting and Merging)

- B1: Chia

- Chia ảnh thành 4 phần
- Kiểm tra tính đơn nhất trong mỗi phần ($\max P - \min P \leq T$)
- Với mỗi vùng không đơn nhất, tiếp tục chia nhỏ

- B2: Hợp

- Từ những vùng đã có, hợp lại thành nhiều vùng lớn hơn dựa vào ngưỡng T

1	1	5	6
2	1	6	7
3	2	7	4
1	0	5	5

Original Image

1	1	5	6
2	1	6	7
3	2	7	4
1	0	5	5

Region splitting into 4 quadrant

1	1	5	6
2	1	6	7
3	2	7	4
1	0	5	5

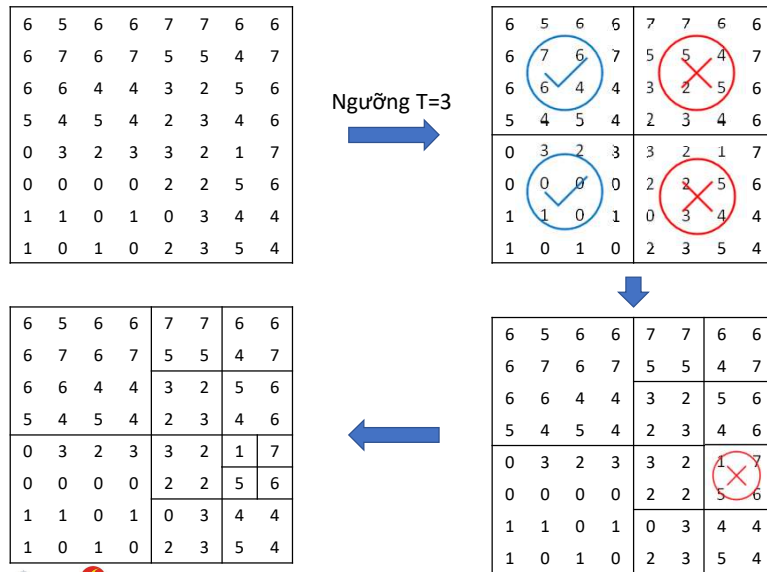
Classifying a quadrant as a region if it satisfies condition else performing further splitting



18

18

Quá trình tách (ví dụ)



VINBIODATA



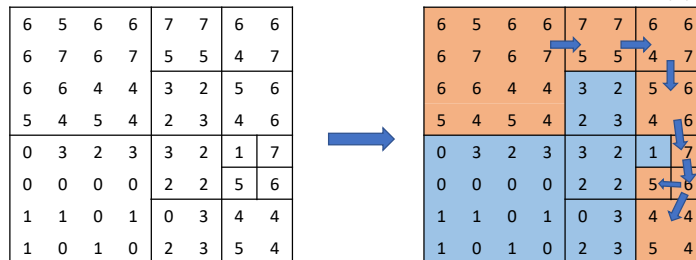
VINGROUP

19

19

Quá trình hợp (ví dụ)

Hợp 2 vùng lân cận nếu
 $\text{Max}(1) - \text{min}(2) \leq T$
 và $\text{max}(2) - \text{min}(1) \leq T$



Ảnh gốc



Ảnh sau khi tách



Ảnh sau hợp



VINBIODATA



VINGROUP

20

20

Các giải thuật phân cụm

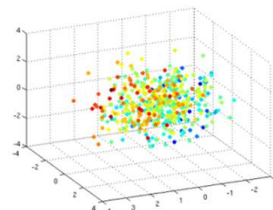


21

21

Segmentation as Clustering

- Điểm ảnh được biểu diễn trong không gian nhiều chiều
 - Màu sắc: 3d
 - Màu sắc + vị trí: 5d



22

22

Giải thuật K-means

- B1: Ngẫu nhiên khởi tạo k điểm trung tâm c_1, c_2, \dots, c_k
- B2: Với các điểm trung tâm hiện tại của từng phân cụm, xác định các điểm còn lại thuộc phân cụm nào
 - Với mỗi điểm p , tìm điểm trung tâm gần nhất và thêm p vào cụm tương ứng
- B3: Với từng phân cụm, tính toán lại điểm trung tâm
 - Điểm trung tâm là trung bình các điểm trong phân cụm
- B4: Nếu c_i thay đổi, lặp tại bước 2

Đặc điểm:

- Luôn hội tụ về một phương án nào đó
- Có thể là “cực tiểu cục bộ”
 - Không đảm bảo luôn tìm thấy hàm mục tiêu cực tiểu toàn cục

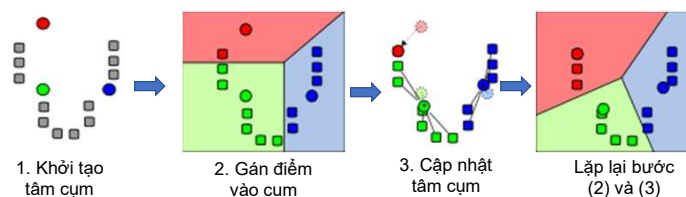
$$\sum_{i=1}^k \sum_{x \in C_i} \|x - c_i\|^2$$



23

23

Minh họa giải thuật K-means



- Demo giải thuật K-means:

<https://www.naftaliharris.com/blog/visualizing-k-means-clustering>

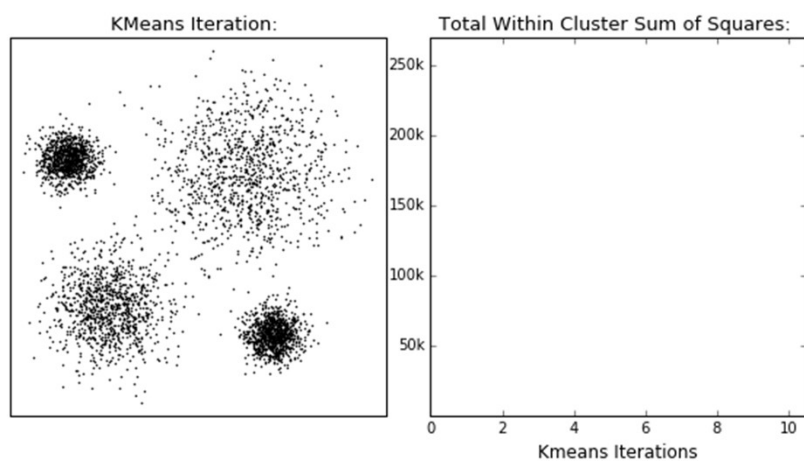
http://home.dei.polimi.it/matteucc/Clustering/tutorial_html/AppletKM.html



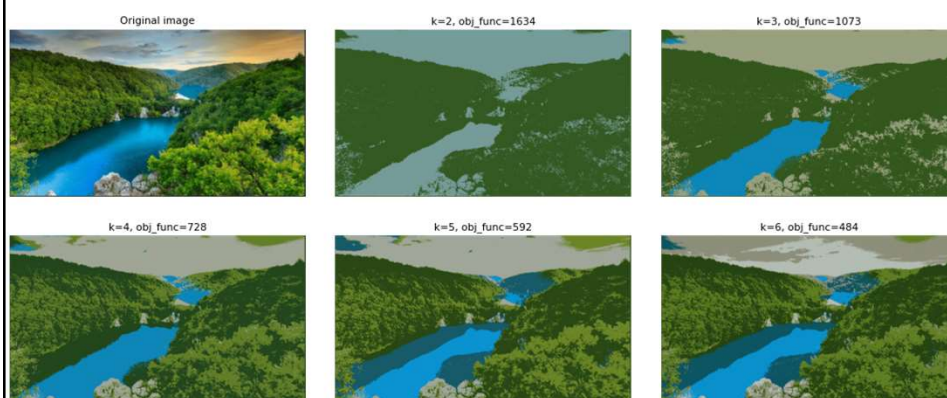
24

24

Minh họa giải thuật K-means (2)

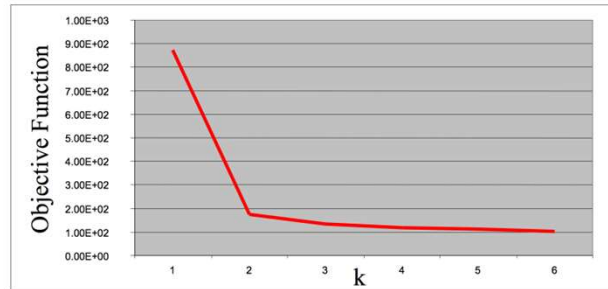


Ví dụ áp dụng K-mean lên ảnh



Lựa chọn số cụm như thế nào?

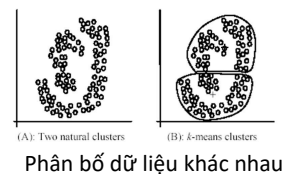
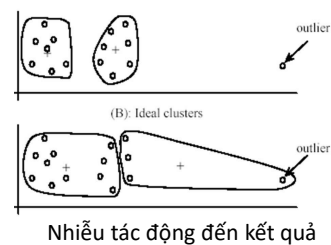
- Thử các giá trị K khác nhau và quan sát kết quả trên tập validation.



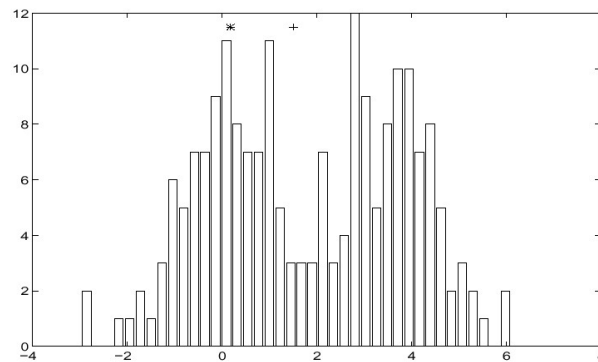
- Có sự thay đổi đột ngột hàm mục tiêu tại $k=2$, nhiều khả năng số nhóm là 2.

K-Means: Ưu và nhược điểm

- Ưu điểm
 - Đơn giản, nhanh, dễ cài đặt
- Nhược điểm
 - Cần chọn K
 - Nhạy cảm với ngoại biên (outliers)
 - Hội tụ về lời giải địa phương
 - Làm việc không tốt nếu phân bố dữ liệu các cụm không giống nhau
 - *Có thể chậm: $O(KNd)$
- Sử dụng
 - Phân cụm không giám sát
 - Ít dùng cho phân đoạn ảnh hoặc là bước trung gian cho phân đoạn



Thuật toán Mean-Shift



- Tìm kiếm đỉnh (mode) bằng thủ tục lặp
 - Khởi tạo các điểm ngẫu nhiên, và cửa sổ W
 - Tính trọng tâm ("mean") của cửa sổ W : $\sum_{x \in W} x H(x)$
 - Tịnh tiến cửa sổ tìm kiếm tới trọng tâm
 - Lặp lại bước 2 cho tới khi hội tụ



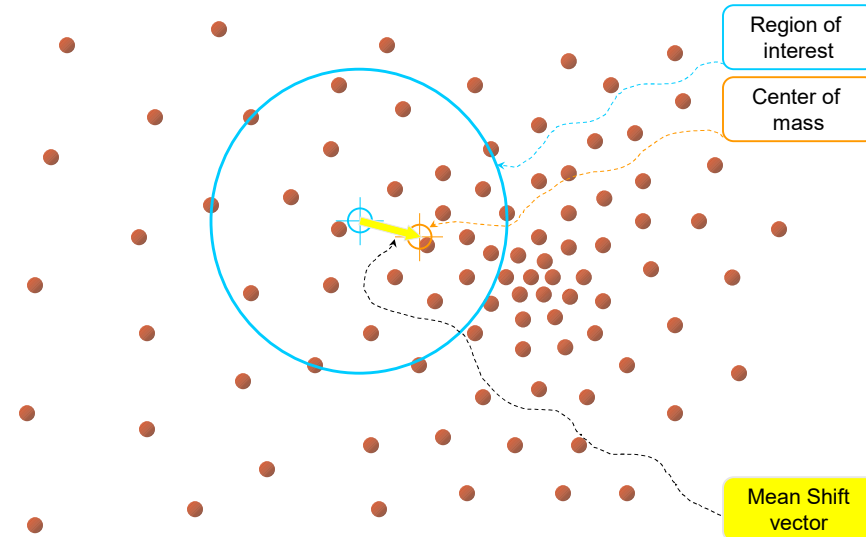
VINBIODATA



29

29

Mean-Shift



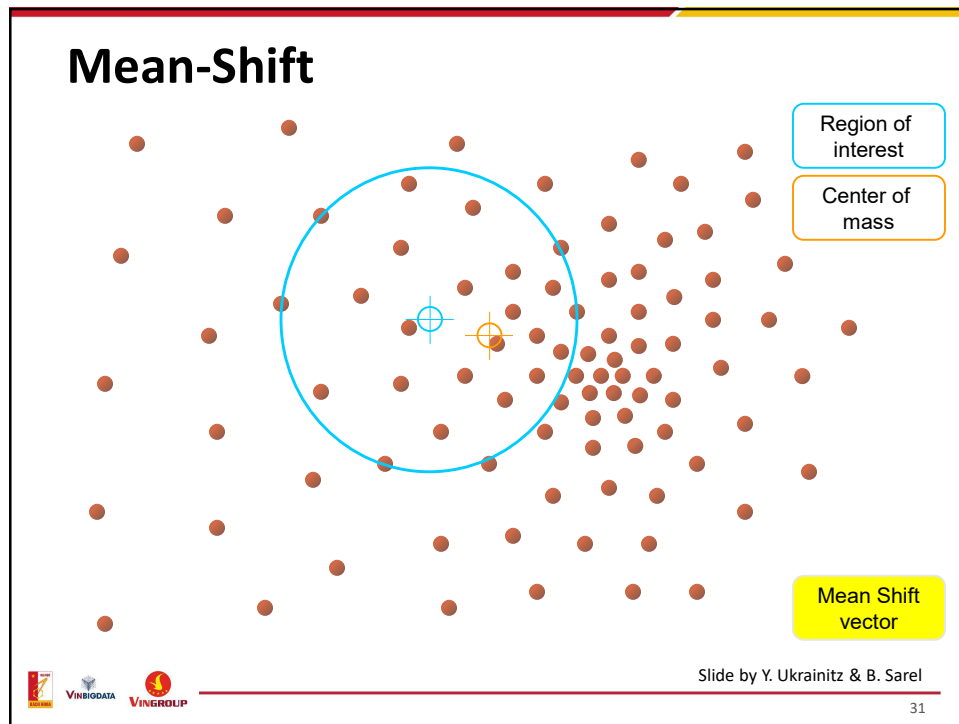
VINBIODATA



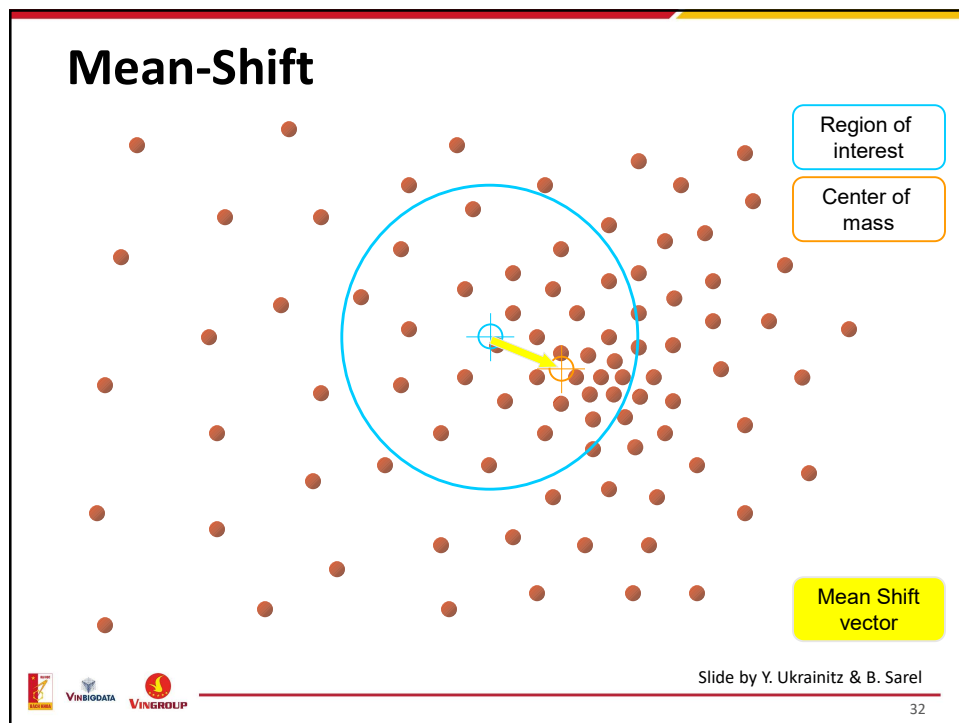
Slide by Y. Ukrainitz & B. Sarel

30

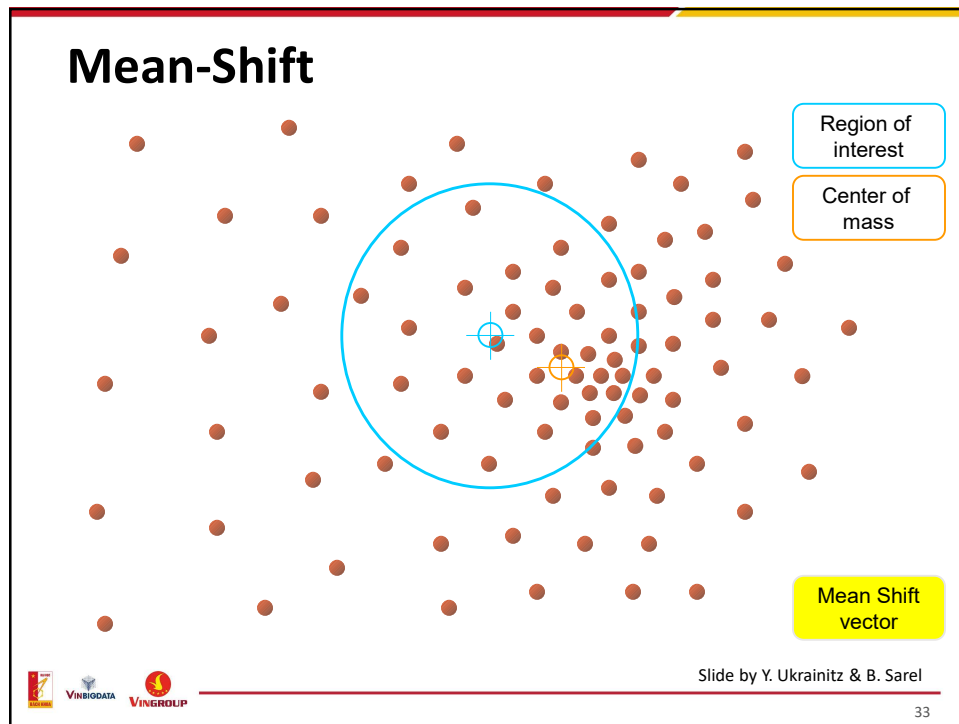
30



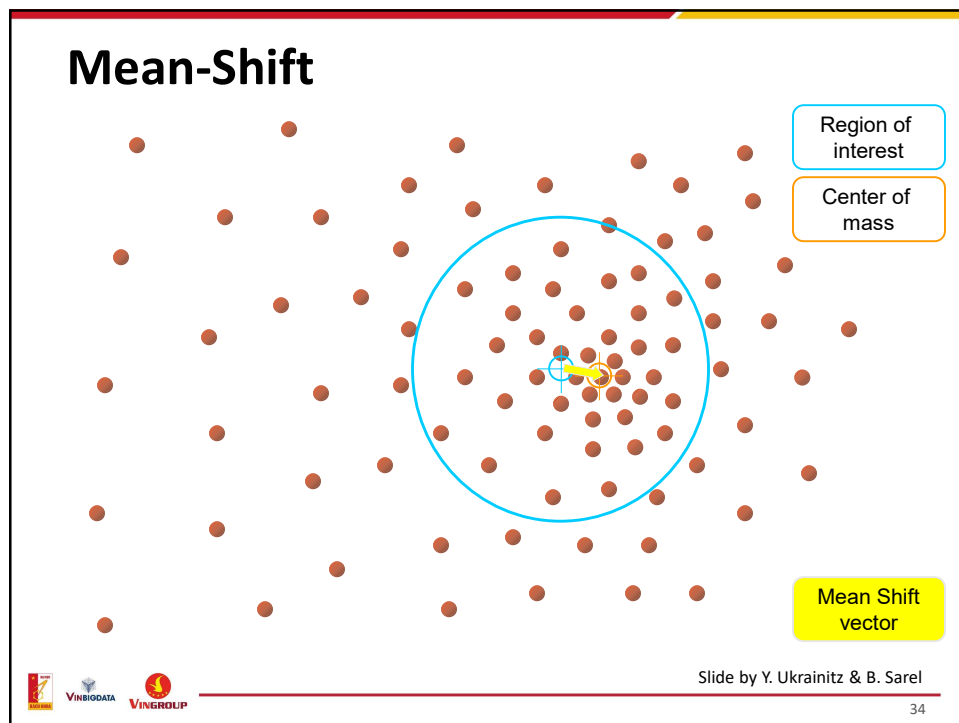
31



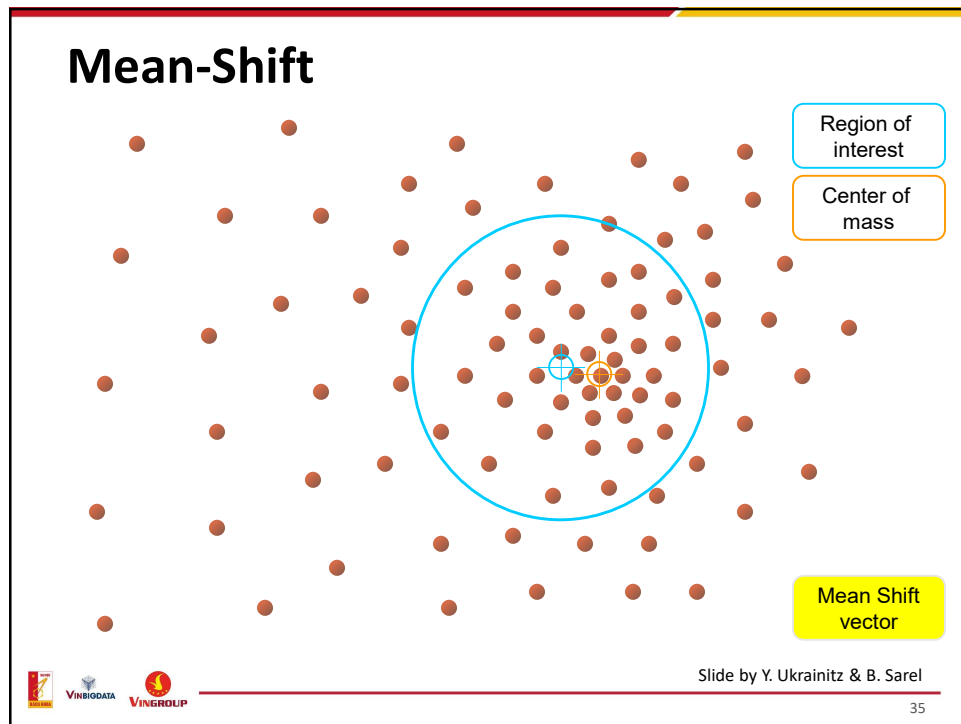
32



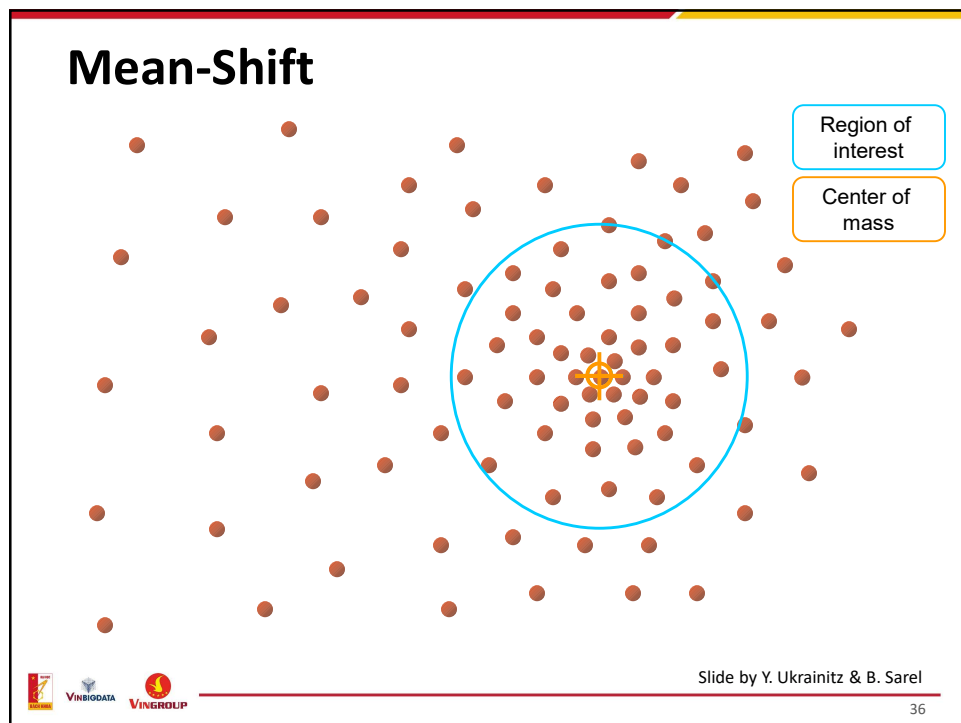
33



34

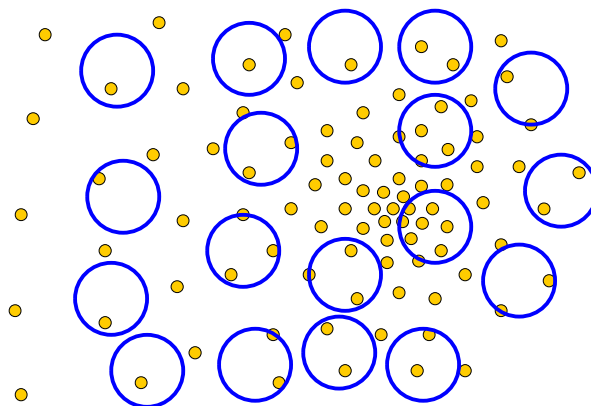


35



36

Real Modality Analysis



Tessellate the space with windows

Chạy song song các quá trình



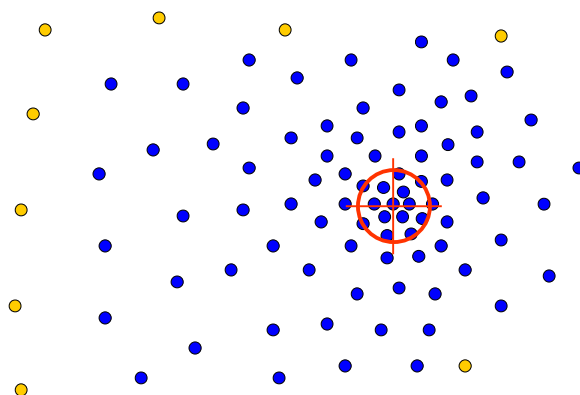
VINBIODATA



37

37

Real Modality Analysis



Các điểm **màu xanh** data là các điểm được các cửa sổ tìm kiếm quét qua trong quá trình hội tụ tới đỉnh (mode).



VINBIODATA

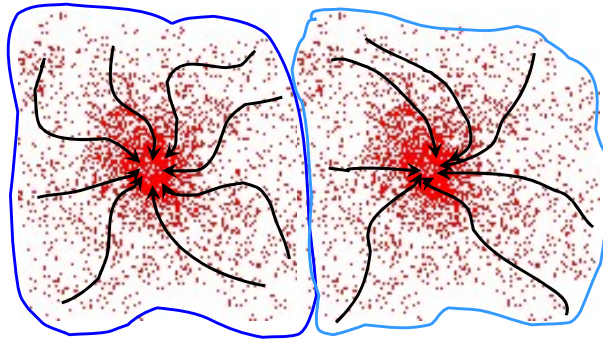


38

38

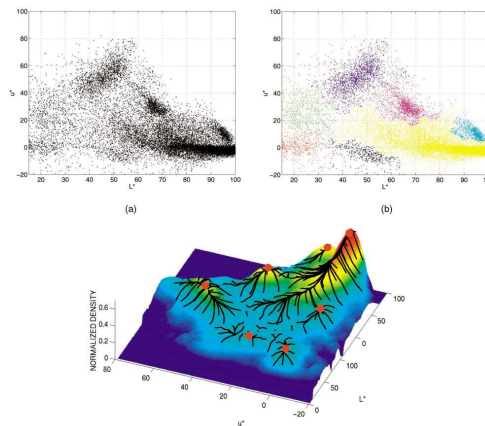
Phân cụm dùng Mean-Shift

- Cluster: tất cả dữ liệu trong vùng “attraction basin” của một đỉnh mode
- Attraction basin: là vùng là tất cả các cửa sổ tìm kiếm trong đó đều hội tụ về cùng một đỉnh.

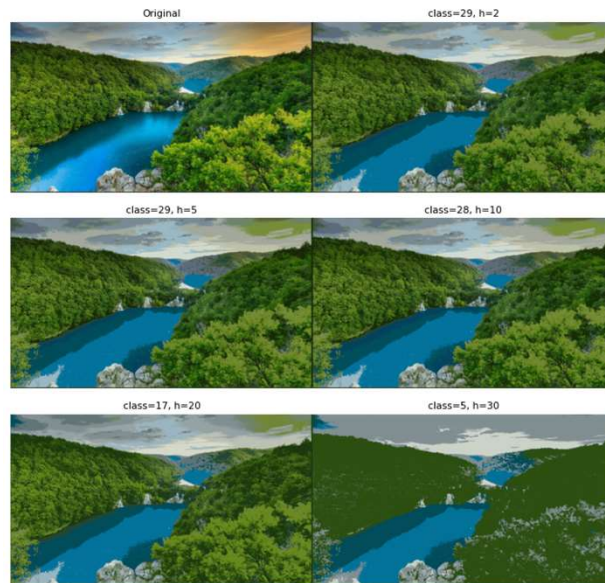


Phân cụm/phân đoạn dùng Mean-Shift

- Trích xuất đặc trưng (color, gradients, texture, etc)
- Khởi tạo các cửa sổ tìm kiếm tại các điểm ảnh khác nhau
- Thực hiện mean shift đối với mỗi cửa sổ cho tới khi hội tụ
- Nhập lại (merge) tất cả các cửa sổ hội tụ về gần cùng một đỉnh



Ví dụ áp dụng Mean-shift lên ảnh



Tổng kết về Mean-Shift

- Ưu điểm
 - Một công cụ tổng quát cho nhiều bài toán
 - Không cần biết phân bố của các cụm dữ liệu (spherical, elliptical, etc.)
 - Chỉ có một tham số duy nhất (kích thước cửa sổ h)
 - Có thể tìm thấy nhiều cụm
 - Không bị ảnh hưởng bởi dữ liệu ngoại lai outliers
- Nhược điểm
 - Kết quả phụ thuộc kích thước cửa sổ
 - Lựa chọn kích thước cửa sổ là không tầm thường
 - Tính toán chậm ($\sim 2s/\text{ảnh}$)
 - Không có khả năng mở rộng tốt khi số chiều không gian đặc trưng tăng

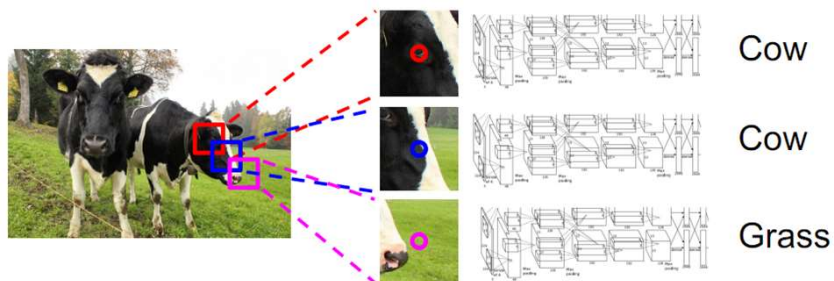
Kỹ thuật học sâu



43

43

Phương pháp cửa sổ trượt



- Phù hợp cho bài toán nhận dạng đối tượng
- Không hiệu quả trong phân đoạn ảnh

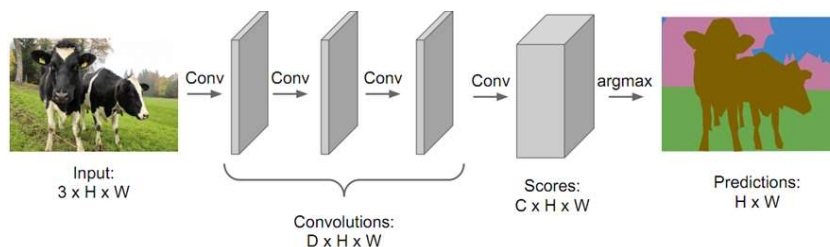


44

44

Tích chập hoàn toàn (FCN)

- Vấn đề: Tích chập ở độ phân giải cao tốn chi phí



45

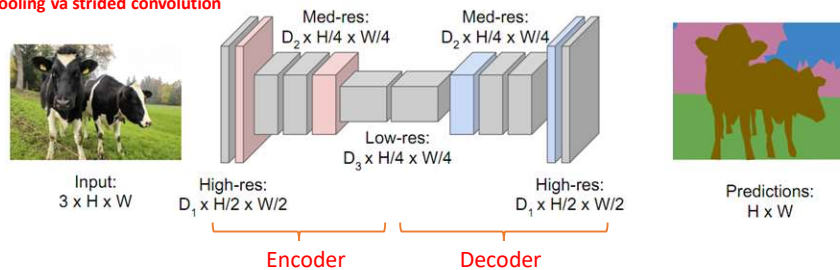
45

Tích chập hoàn toàn (FCN)

- Thiết kế mạng có nhánh downsampling và upsampling

Downsampling:
pooling và strided convolution

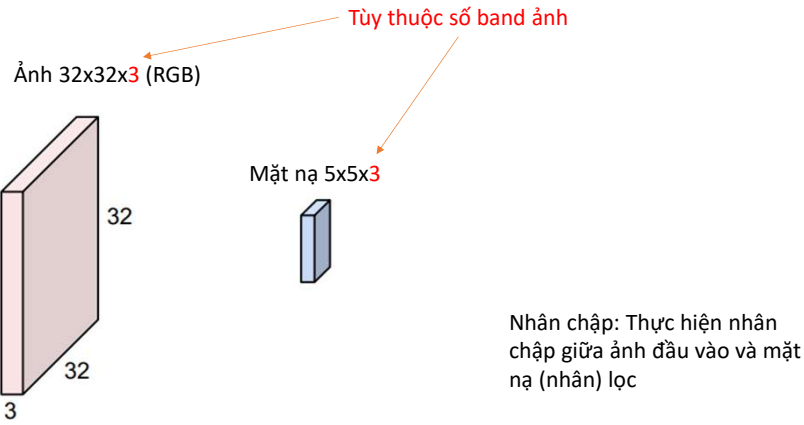
Upsampling:



46

46

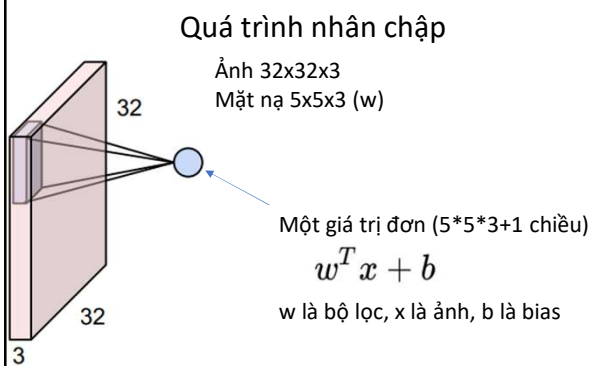
Lớp tích chập



47

47

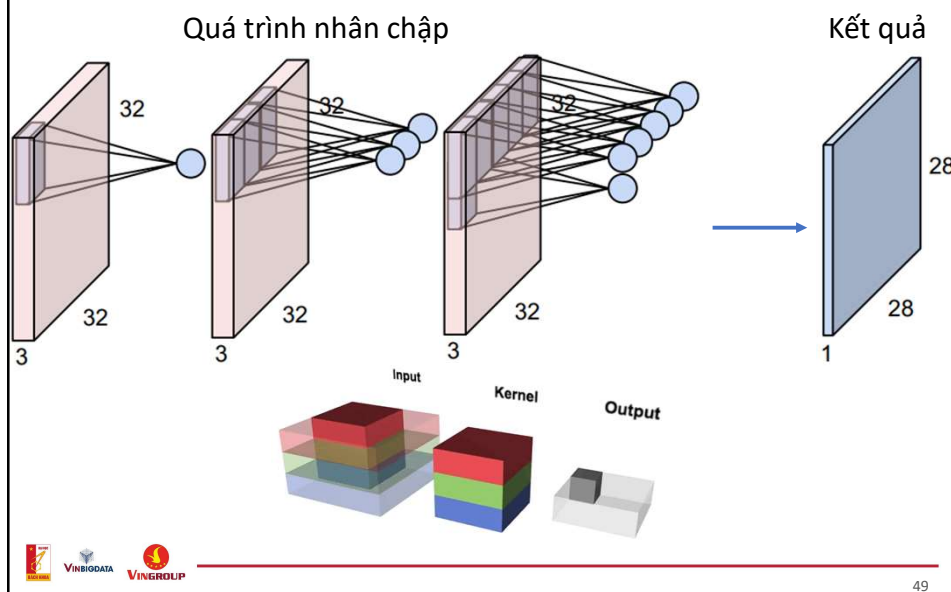
Lớp tích chập



48

48

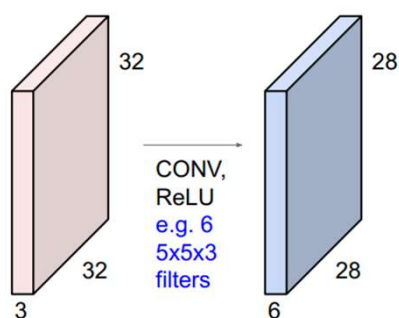
Lớp tích chập



49

Lớp tích chập

- Kết hợp với các hàm kích hoạt (Activation function)



50

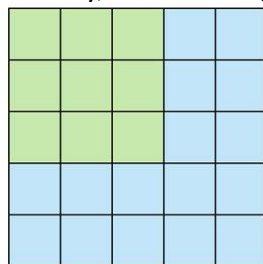
Padding và Stride

- Padding

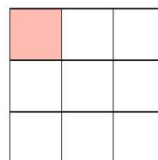
- Thêm viền bên ngoài để tăng kích thước ảnh đầu ra
- Thường dùng để giữ kích thước ảnh đầu ra giống đầu vào

- Stride

- Bước nhảy, bình thường là 1, có thể lớn hơn 1



Stride 1



Feature Map



VINBIODATA



Stride 1 và pad 0

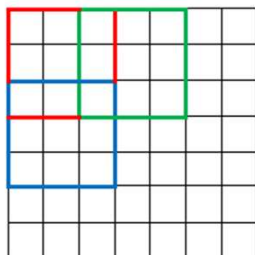
51

51

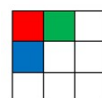
Padding và Stride (2)

Stride = (2, 2)

Ảnh đầu vào 7x7



Ảnh đầu ra 3x3



VINBIODATA

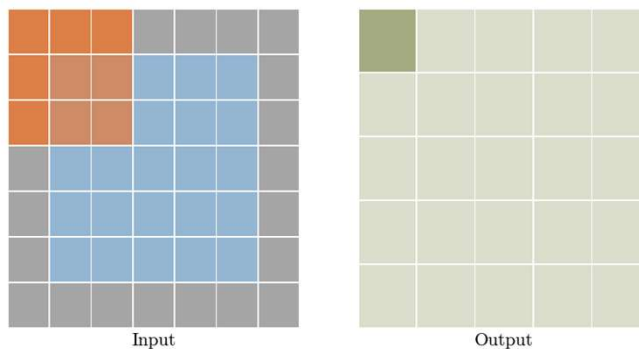


52

52

Padding và Stride (2)

Type: conv - Stride: 1 Padding: 1



Padding: zero padding và nearest neighbour

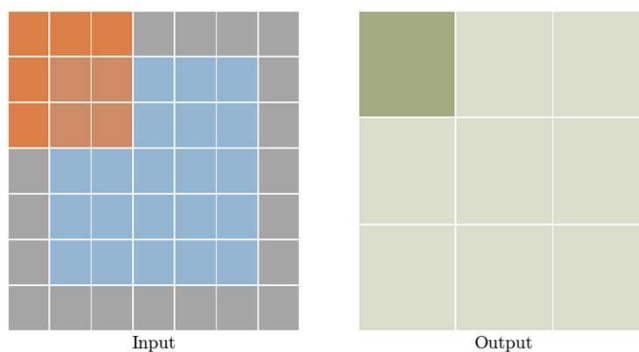


53

53

Padding và Stride (4)

Type: conv - Stride: 2 Padding: 1



55

55

Pooling layer

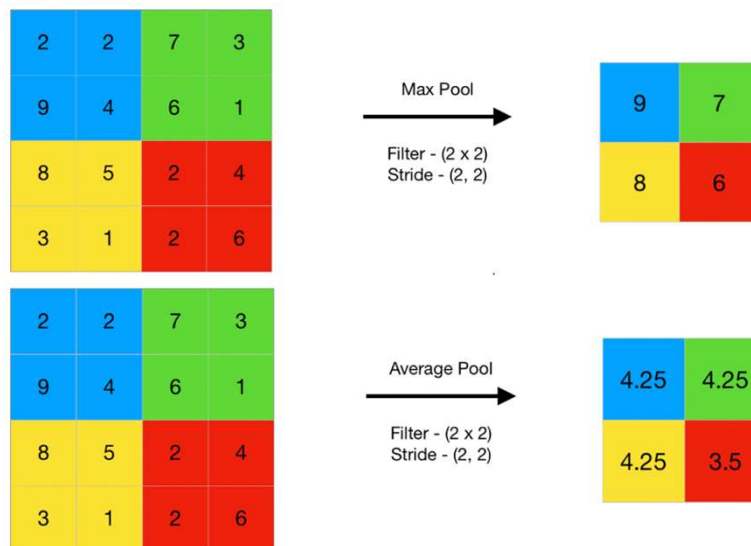
- Là kỹ thuật để giảm kích thước ảnh
- Các kiểu pooling
 - Max (cực đại): giá trị lớn nhất trong cửa sổ
 - Average (trung bình): giá trị trung bình trong cửa sổ
 - Global (toàn cục): mỗi kênh chuyển về một giá trị duy nhất (cực đại hoặc trung bình)



56

56

Pooling layer (ví dụ)



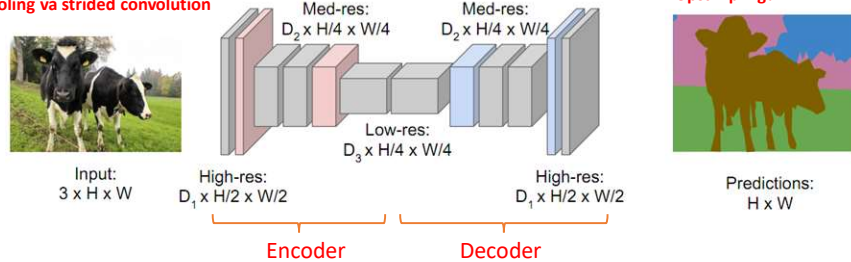
57

57

Tích chập hoàn toàn (FCN)

- Thiết kế mạng có nhánh downsampling và upsampling

Downsampling:
pooling và strided convolution



58

58

Unpooling

- Thực hiện ngược lại pooling

Nearest Neighbor

1	2
3	4

Input: 2×2

1	1	2	2
1	1	2	2
3	3	4	4
3	3	4	4

Output: 4×4

"Bed of Nails"

1	2
3	4

Input: 2×2

1	0	2	0
0	0	0	0
3	0	4	0
0	0	0	0

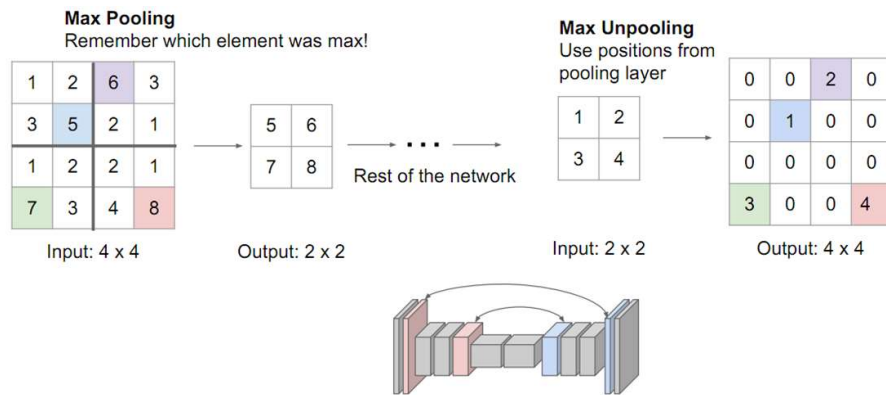
Output: 4×4



59

59

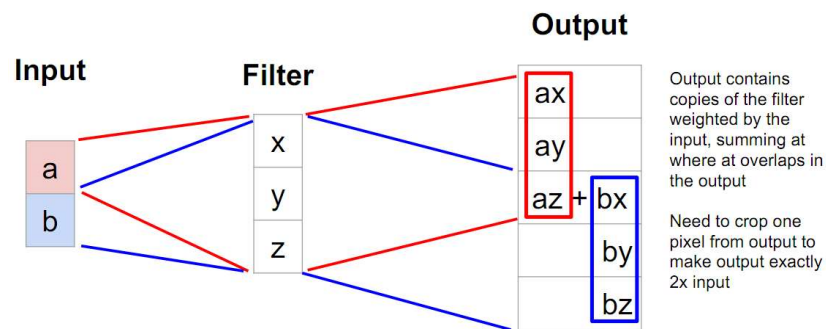
“Max Unpooling”



Transpose convolution layer

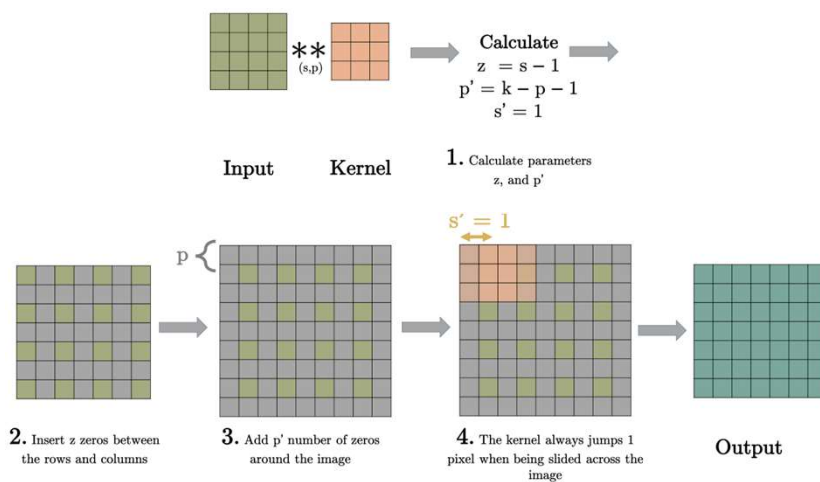
- Thực hiện ngược lại phép nhân chập để nâng độ phân giải.
- Một số tên gọi khác
 - Deconvolution (không nên dùng)
 - Upconvolution
 - Backward strided convolution
 - Fractionally strided convolution

Transpose convolution: Ví dụ 1D



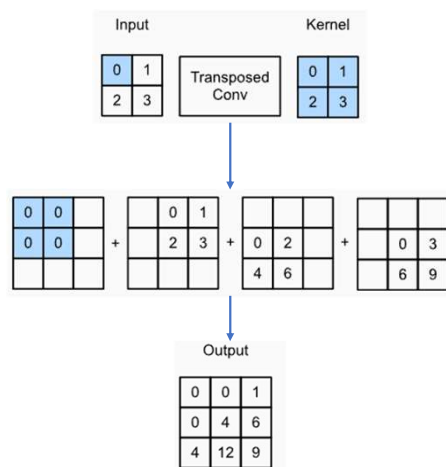
62

Tranpose convolution: 2D

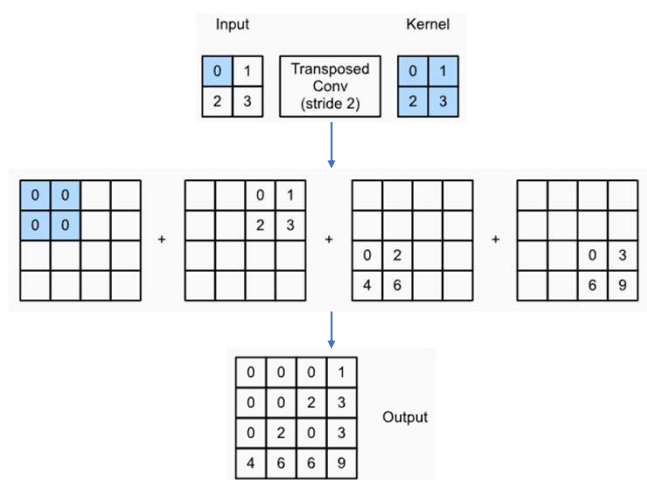


63

Tranpose convolution: ví dụ 2D

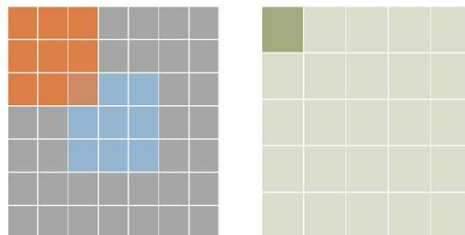


Tranpose convolution: ví dụ 2D

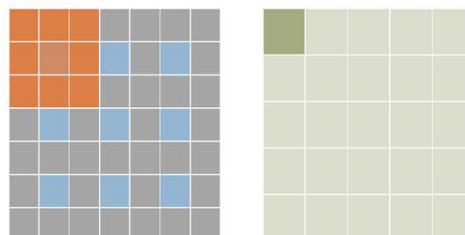


Tranpose convolution: ví dụ

Type: transposed conv - Stride: 1 Padding: 0



Type: transposed conv - Stride: 2 Padding: 1



Input

Output



66

66

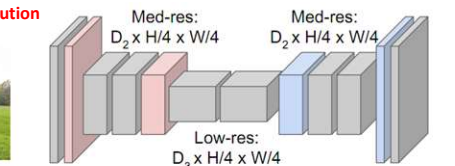
Tích chập hoàn toàn (FCN)

- Thiết kế mạng có nhánh downsampling và upsampling

Downsampling:
pooling và strided convolution



Input:
 $3 \times H \times W$



Encoder

Decoder

Upsampling:
Unpooling và strided
transpose convolution



Predictions:
 $H \times W$

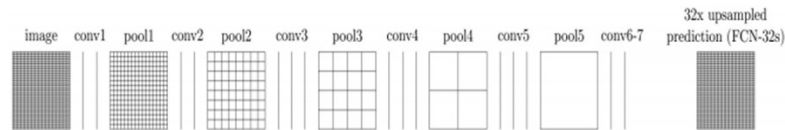


67

67

Skip connection (kết nối tắt)

- Quá trình downsampling/ mã hóa làm giảm độ phân giải đầu vào → quá trình upsampling/ giải mã khó phân đoạn chi tiết ảnh



Ảnh đầu vào

Ground truth

Kết quả phân đoạn

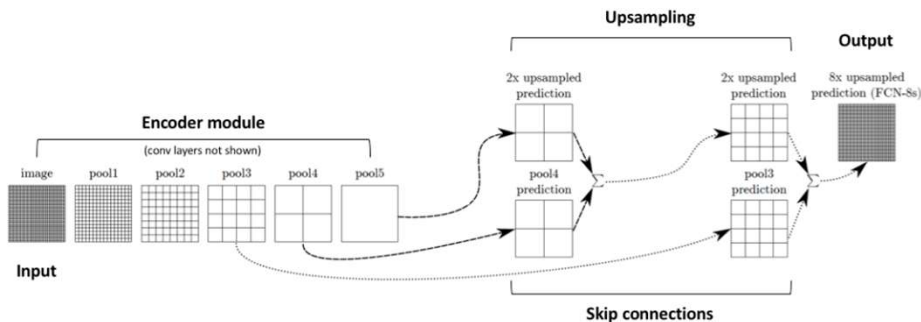


68

68

Skip connection (kết nối tắt)

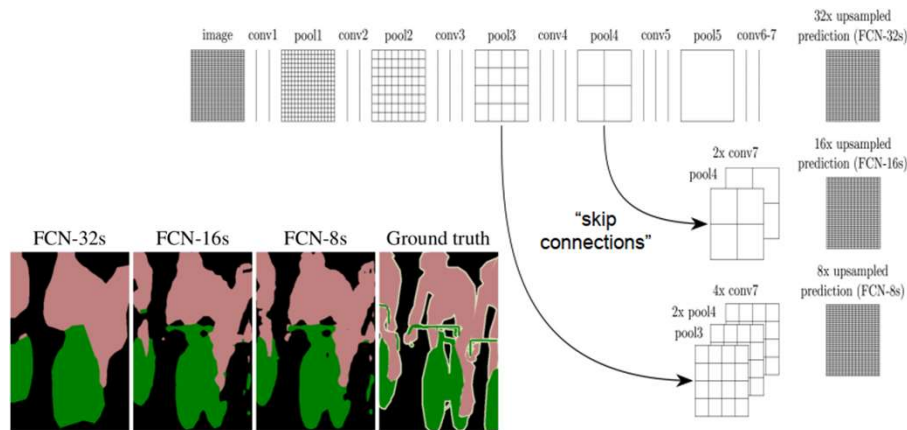
- Thêm các kết nối tắt các dữ liệu độ phân giải cao hơn (đã bị loại bỏ) vào kết quả trong quá trình upsampling → tăng độ chi tiết phân loại



69

69

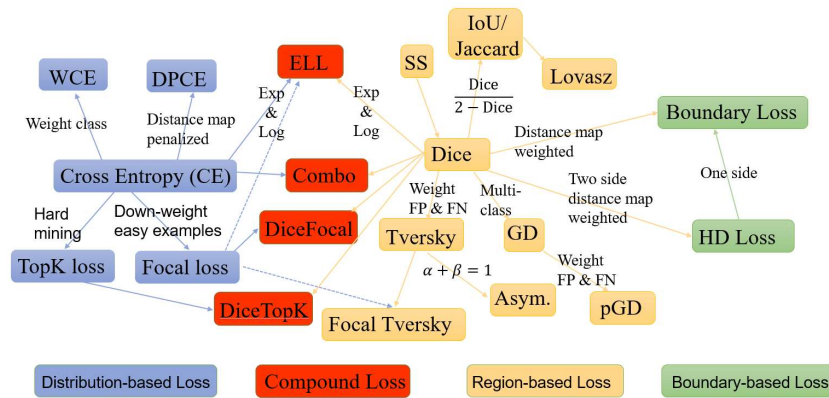
Kết quả giữa có và không có kết nối tắt



Lưu ý khi dùng kết nối tắt

- Không gian giữ liệu phải giống nhau
- Hai phép toán phổ biến
 - Cộng (addition)
 - Nối (concatenation)
- Hai loại kết nối tắt
 - Kết nối tắt ngắn (short skip connection) : những lớp conv liên tiếp không có sự thay đổi kích thước dữ liệu (vd: ResNet)
 - Kết nối tắt dài (long skip connection): dùng cho những mạng đối xứng có sự thay đổi kích thước dữ liệu, cần áp dụng cho từng lớp mã hóa/ giải mã tương ứng cùng kích thước (FCN, Unet)

Hàm mục tiêu



72

72

Hàm mục tiêu dựa trên phân phối

- Cross Entropy (CE):

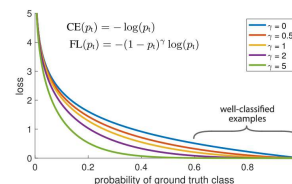
$$CE(p, \hat{p}) = -(p \log(\hat{p}) + (1 - p) \log(1 - \hat{p}))$$

- Weighted CE: mỗi lớp có trọng số khác nhau

$$WCE(p, \hat{p}) = -(\beta p \log(\hat{p}) + (1 - p) \log(1 - \hat{p}))$$

- Focal loss: giải quyết vấn đề mất cân bằng lớn giữa lớp nền và lớp đối tượng quan tâm. Giá trị hàm mục tiêu đối với những mẫu dễ phân loại được giảm xuống thấp để mạng tập trung hơn vào mẫu khó.

$$FL(p, \hat{p}) = -(\alpha(1 - \hat{p})^\gamma p \log(\hat{p}) + (1 - \alpha)\hat{p}^\gamma(1 - p) \log(1 - \hat{p}))$$



73

73

Hàm mục tiêu dựa trên vùng

- Dice coefficient và IoU:

$$DC = \frac{2TP}{2TP + FP + FN} = \frac{2|X \cap Y|}{|X| + |Y|}$$

$$IoU = \frac{TP}{TP + FP + FN} = \frac{|X \cap Y|}{|X| + |Y| - |X \cap Y|}$$

- Dice loss: $DL(p, \hat{p}) = 1 - \frac{2p\hat{p} + 1}{p + \hat{p} + 1}$

- Tversky loss:

$$TI(p, \hat{p}) = \frac{p\hat{p}}{p\hat{p} + \beta(1-p)\hat{p} + (1-\beta)p(1-\hat{p})}$$



74

74

Hàm mục tiêu kết hợp

- Dice loss + CE:

$$CE(p, \hat{p}) + DL(p, \hat{p})$$

- Dice loss + Focal loss

$$CE(p, \hat{p}) + FL(p, \hat{p})$$

- ...



75

75

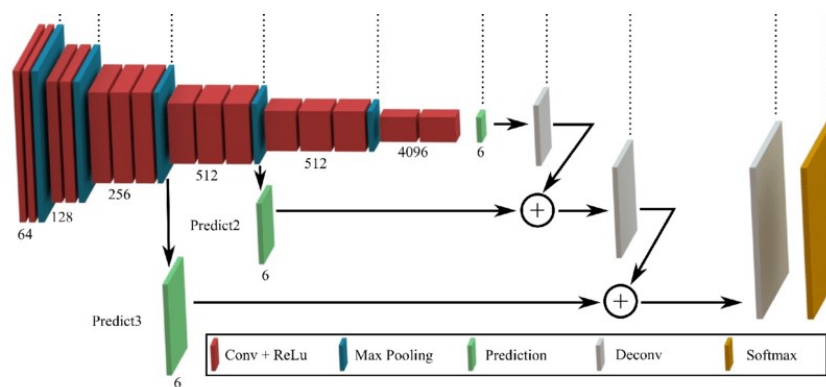
Một số mạng phân đoạn ảnh tiêu biểu



78

78

FCN với 2 kết nối tắt

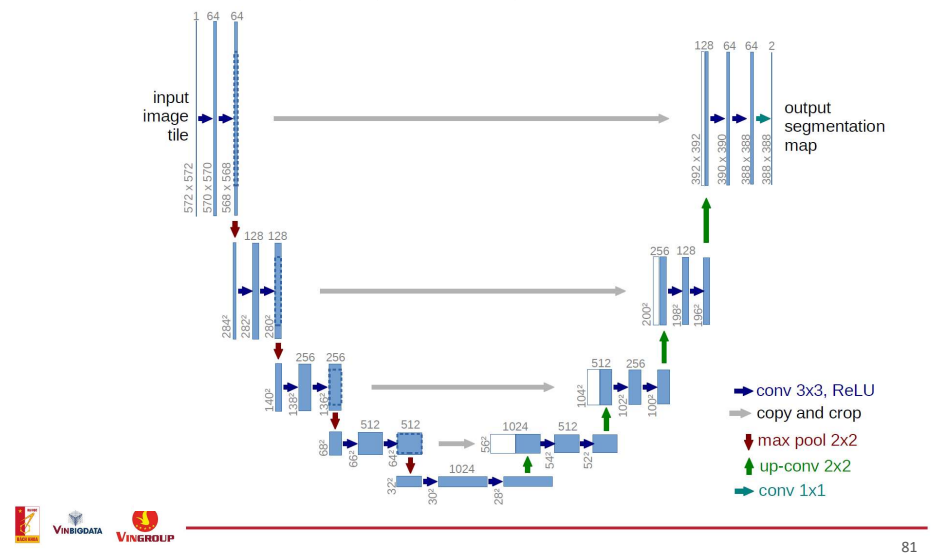


79

79

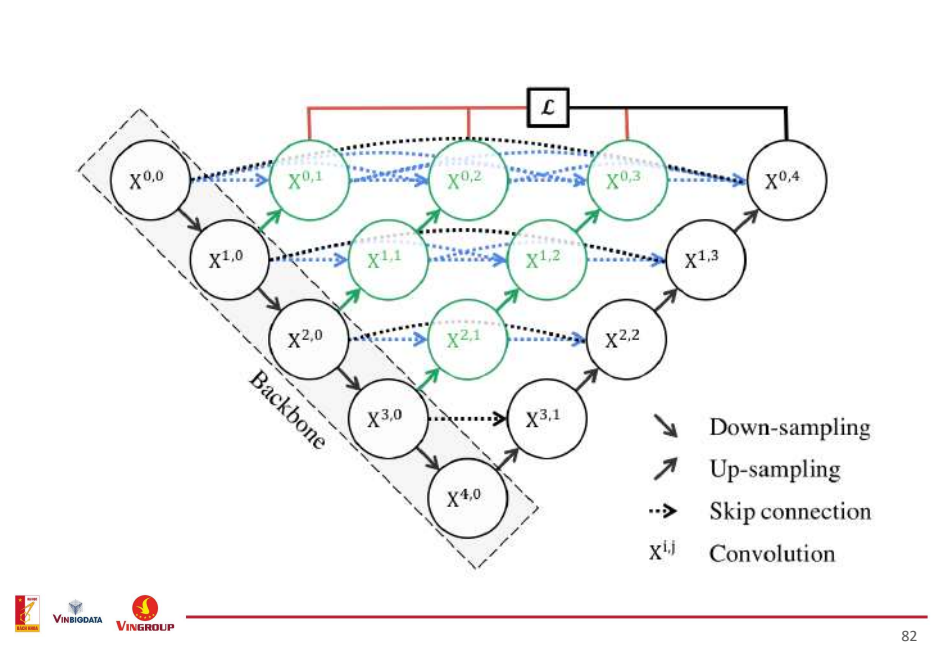
U-Net

- Được sử dụng rộng rãi trong y tế



81

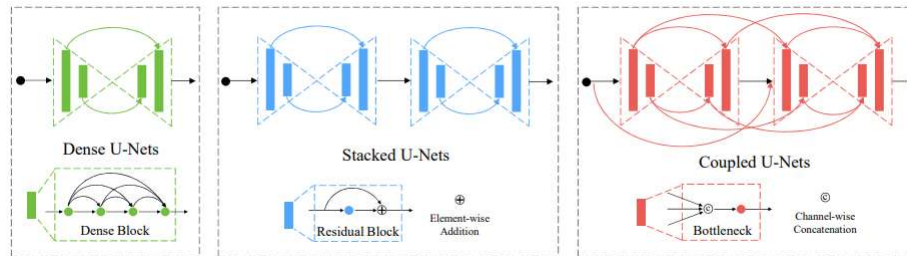
U-Net++



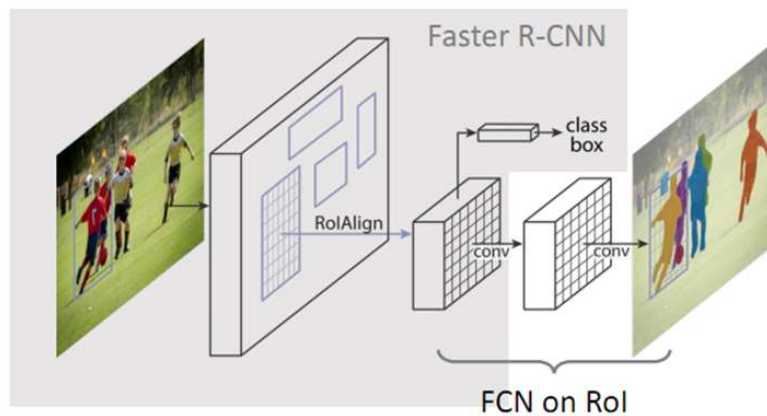
82

Stacked UNets và CUNets

- Stacked UNets: ghép nhiều UNet nối tiếp nhau
- CUNets: cũng ghép nhiều UNet nối tiếp nhau nhưng có thêm các kết nối tắt giữa các UNet với nhau



Mask R-CNN



Transformer-based models

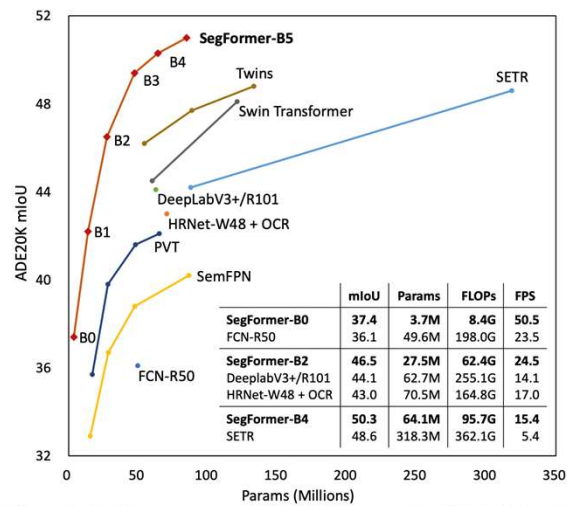


Figure 1: **Performance vs. model efficiency on ADE20K.** All results are reported with single model and single-scale inference. SegFormer achieves a new state-of-the-art 51.0% mIoU while being significantly more efficient than previous methods.



85

85

Tài liệu tham khảo

- Kristen Grauman (CS 376: Computer Vision, Spring 2018, The University of Texas at Austin)
- Stanford CS231n Course:
<http://cs231n.stanford.edu/>



87

87

