

Q1) a)  $H \times W \times M$  kernel  $k \times k \times M$   $\rightarrow N$

same padding - stride = 1  $\Rightarrow N \times (k \times k \times M + 1) =$  مقدار پیکسلها

تعداد ضربهای  
حسابی  $\rightarrow k \times k \times M$

تعداد جمع  $\rightarrow k \times k \times M - 1$

تعداد اعلان شدن کرنل  $\rightarrow H \times W \times N \Rightarrow$  ضرب  $H \times W \times N \times k^2 \times M$   
جمع  $\rightarrow H \times W \times N \times (k^2 \times M - 1)$

b) input size  $\rightarrow 128 \times 128 \times 3$  kernel:  $5 \times 5$

تعداد پیکسلها  $\rightarrow 64 \times (5 \times 5 \times 3 + 1) = 4964$   $\rightarrow$  مقدار پیکسلها

$$H_{out} = \left\lfloor \frac{H_{in} + 2p - k}{s} + 1 \right\rfloor = \left\lfloor \frac{128 + 2 \times 2 - 5}{2} + 1 \right\rfloor = 64$$

$\Rightarrow$  ابعاد خروجی  $= 64 \times 64 \times 64$   $\Rightarrow$  مقدار ضرب  $= 64 \times 64 \times 64 \times 25 \times 3 = 19660800$   
تعداد جمع  $= 64 \times 64 \times 64 \times (25 \times 3 - 1) = 19378656$

تعداد پیکسلها:  $128 \times (5 \times 5 \times 64 + 1) = 204928 \rightarrow$  مقدار پیکسلها

ابعاد:  $\left\lfloor \frac{64 + 4 - 5}{2} + 1 \right\rfloor = 32 \Rightarrow$  output size  $= 32 \times 32 \times 128$

ضرب  $= 32 \times 32 \times 128 \times 5^2 \times 64 = 209715200$

جمع  $= 32 \times 32 \times 128 \times (5^2 \times 64 - 1) = 209584128$

تعداد پیکسلها:  $256 (5 \times 5 \times 128 + 1) = 819456$

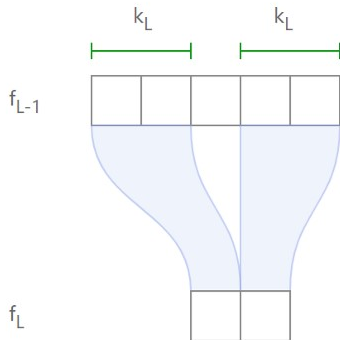
$\left\lfloor \frac{32 + 4 - 5}{2} + 1 \right\rfloor = 16 \Rightarrow$  output size  $= 16 \times 16 \times 256$

$\Rightarrow$  ضرب  $= 16 \times 16 \times 256 \times 5^2 \times 128 = 209715200$

جمع  $= 16 \times 16 \times 256 \times (5^2 \times 128 - 1) = 209649664$

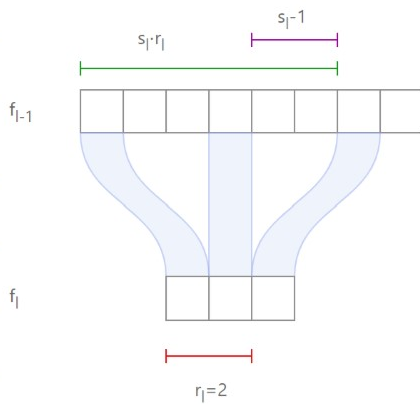
receptive field :  $\Rightarrow r_L = 1$  تعداد منفرجه‌های از لایه  $L$  که در شکل سبز  $\rightarrow r_L$  یک منفرجه آخری دارند

نوع ایجاد شدن هر منفرجه:



Kernel Size ( $k_L$ ): 2  
Padding ( $p_L, q_L$ ): 0  
Stride ( $s_L$ ): 3

فرض است  $k_L = 1$ :

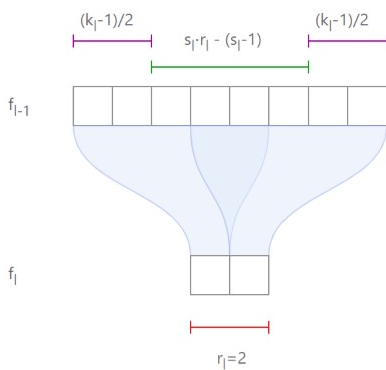


Kernel Size ( $k_L$ ): 1  
Padding ( $p_L, q_L$ ): 0  
Stride ( $s_L$ ): 3

$$r_{L-1} = s_L r_L - (s_L - 1)$$

حالا بفرمایید از شکل مشخص است.

$k_L < 1$ :



Kernel Size ( $k_L$ ): 5  
Padding ( $p_L, q_L$ ): 0  
Stride ( $s_L$ ): 3

$$r_{L-1} = s_L r_L + (k_L - s_L)$$

لازم  $\Rightarrow$  به صورت بازرسی  $r_0 = \sum_{l=1}^L (k_l - 1) \prod_{i=1}^{l-1} s_i + 1$

$L = 3 \rightarrow$  برای سوال  $S = 2$  و  $k = 5$  برای لایه‌ها

$$\Rightarrow r_0 = 1 + 1 \times 2 + 1 \times 2 \times 2 + 1 = 29$$

c) 1) Depthwise  $\rightarrow$  bottleneck  $k \times k$  و  $1 \times 1 \times M$  (بسیار)

input size:  $H \times W \times M$

$\Rightarrow M \times (k^2 + 1) \rightarrow$  مرحله اول  $(1 \times 1 \times M + 1) \rightarrow$  مرحله دوم

$\Rightarrow M \times (k^2 + 1) + M + 1 = M(k^2 + 2) + 1$

پارامترها  $\rightarrow$  کمتر نسبت به  $3D$   $\rightarrow$  ضرب  $M \times W \times M \times k^2$

2)

kernel:  $5 \times 5$  دروس:  $64 \times 64 \times 64$  : لایه دوم

$\Rightarrow$  پارامتر  $64(25 + 2) + 1 = 1729$

$\left\lfloor \frac{64 + 4 - 5}{2} + 1 \right\rfloor = 32$  output  $\rightarrow 32 \times 32$

$\Rightarrow$  ضرب  $32 \times 32 \times 64 \times 25 = 1638400$

لایه سوم:  $32 \times 32 \rightarrow$  دروس

پارامتر  $1 \times (25 + 2) + 1 = 28$

$\left\lfloor \frac{32 + 4 - 5}{2} + 1 \right\rfloor = 16$  output  $\rightarrow 16 \times 16$

ضرب  $16 \times 16 \times 1 \times 25 = 6400$

پارامترها و عملیات  $\rightarrow$  در حالتی که توضیحی کمتر است.

تعداد پارامترها  $= 1029248$   
ماندوبش عدد

تعداد پارامترها  $= 6621$   
Depthwise

۱) کانولوشن عادی  $\rightarrow$  output size =  $16 \times 16 \times 256$

Flatten  $\rightarrow$  output size = 65536      200 class output

$$\Rightarrow \text{پارامترهای FC} = 65536 \times 200 + 200 = 13107400$$

$\rightarrow$  total parameter number = 14136648  $\rightarrow$  93.6%  
پارامترهای FC است

Depth wise  $\rightarrow$  output size =  $16 \times 16$

$$\Rightarrow \text{FC پارامترها} = 16 \times 16 \times 200 + 200 = 51400$$

$$\text{کل پارامترها} = 51400 + 6621 = 58021 \rightarrow 88.6\% \text{ سهم پارامترها FC}$$

برای کاهش پارامترها FC، global average pooling، dropout

همان bottleneck layer = Depth wise separable

مزایای این لایه: کاهش پارامترها و محاسبه



(1) نحوه ترکیب ویژگی‌ها: DenseNet ← خروجی هر لایه با خروجی لایه‌های قبلی

Concat ← دستور و در ResNet ← خروجی هر لایه بالا به لایه‌های قبلی جمع می‌شود

استفاده از Concat ← استفاده مجدد از ویژگی‌های را مقوی می‌کند. و جمع ← جریان‌های را با هم  
ماده می‌کند

(2) جریان مستقیم ترابیان: از آنجایی که هر لایه با قبلی ها Concat می‌شود، این تقنین وجود دارد که ترابیان به طور مستقیم صنفی می‌شود و اصالت معنای آن کمتر است.

استفاده مجدد از ویژگی‌ها: از خروجی لایه‌های قبلی استفاده می‌شود که باعث redundancy

و به صرفه‌جویی ترابیان کمک می‌کند. چون بهر از صفی‌ها استفاده می‌کنیم اصالت به لایه‌های

عمیق کمک شده و gradient vanishing کم می‌شود

(3) داده‌های کم ← با استفاده مجدد از ویژگی‌ها overfit را در مقایسه با سایر مدل‌ها کم می‌کند

و این به روش‌هایی به طور موثرتری بجهت می‌برد.

عملکرد بهتر با پارامترهای کمتر دارد ← پس اگر محدودیت computational داشته باشیم، استفاده کنیم

چون دقت خوبی دارد.

در image classification و object detection استفاده می‌شود.

مدل کاربردی: تشخیص ناهنجاری‌ها از تصاویر پزشکی ← X-ray ، MRI ، CT

۱۴ برای هر نوع داده branch متفاوتی از DenseNet استفاده می‌کنیم و فیلترها را

استخراج می‌کنیم بعد در یک بلاک مخصوص به Fusion block می‌آمیزیم با concatenate کردن فیلترها

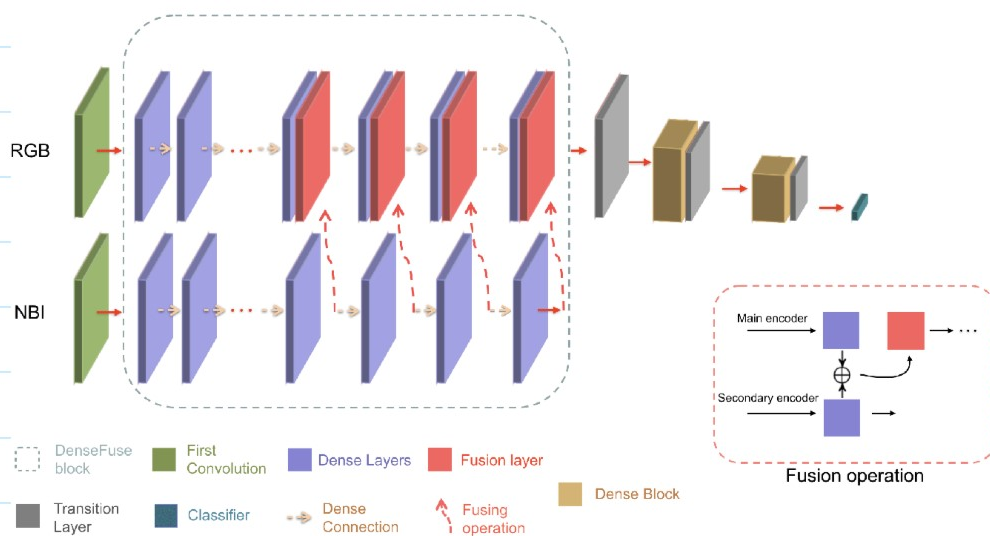
آنها را ترکیب کرده و سپس با سایر لایه‌های DenseNet، به سبب خاصی را افزایش می‌دهیم.

⋮

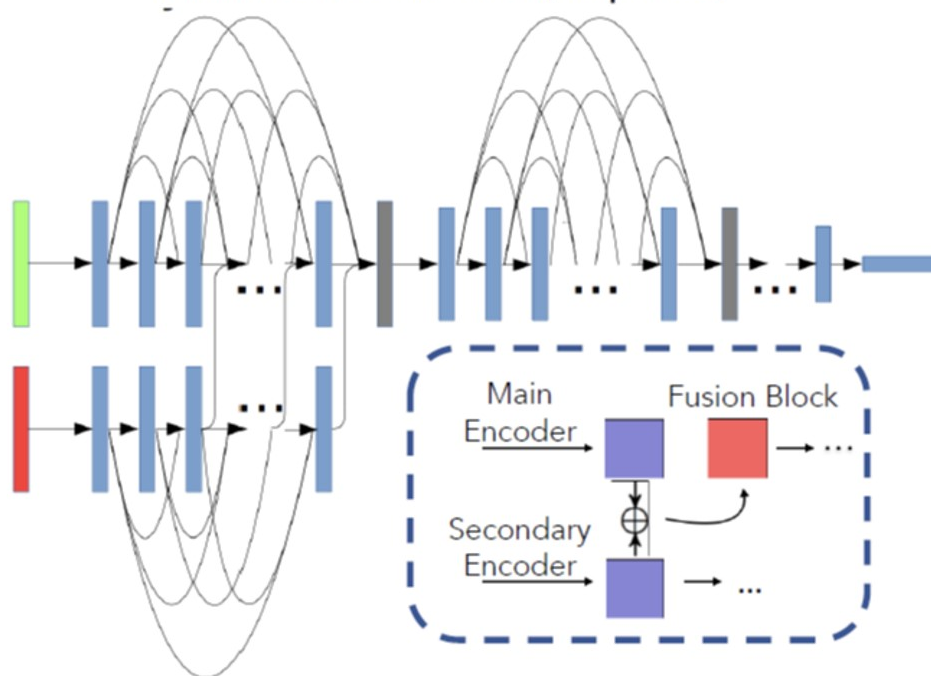
Multi modality input → different Densenet branches → Fusion block → Dense layer → output

⋮

نمودار:



Multimodal Densenet (Proposed)



Q3)

۱۴) با احتمال بزرگ میگوییم از encoder مستقیماً به decoder باعث می‌شود می‌گوییم

قطع می‌کنند (خیز ناگهانی) و به صورت decoding و طبقه بندی می‌شود

با این داده‌ها احتمال است که به اصطلاح می‌زنیم جریان گرادیان به بعد دانه  $\text{gradient vanishing}$

با احتمال کمتری رخ می‌دهد. همچنین باعث می‌شود که آموزش موثرتر انجام شده و همگامی

سریع‌تر اتفاق بیفتد. اطلاعات contextual و spatial را ترکیب می‌کنند و

به این ترتیب باعث می‌شود segmentation انجام می‌شود.

ب) نوعی  $\text{data augmentation}$  است.

چگونه اجرا می‌شود؟ این است که یک سری تغییرات اعمال می‌دهد:

- بردارهای جانبی که یک گرید از تصویر اعمال می‌شود. (مثلاً  $3 \times 3$ )

این جانبی‌ها به صورت تصادفی و با توزیع گاوسی با ابعاد مشخص است.

- جابجایی‌ها با استفاده از Bicubic interpolation به یک نقطه تبدیل می‌شوند تا هر پیکسل

به صورت خاص جابجایی شود.

- این تغییرات بر روی تصویر استفاده می‌شود تا تغییرات «نشتی» تغییر یافته از تغییرات طبیعی مثل

گشودگی و فشردگی و ... ایجاد شود.

تصادفی بودن تغییرات باعث می‌شود واقع‌ترین تجربه باشد.

تاثير کاهش  $\text{over-fit}$  به با  $\text{train}$  شدن روی مجموعه بزرگتر و متنوع‌تری از داده‌ها

$\text{robust}$  شدن به حساسیت به نویز و  $\text{spatial distortion}$  کاهش شده و به تغییرات از این

قبل حساس افراد بود.

اگر  $\text{generalization}$  ← عملکرد بهتر روی (داده‌ها)  $\text{unseen}$

چند داده‌ها را به یاد می‌گیرد که داده‌ها را به  $\text{augmentation}$  بسیار ضروری است

تا بتوانیم با داده‌های محدود عملکرد خوبی داشته باشیم (مثلاً در صورت دیدن یک)

اتریش دمت ← با دیدن داده‌های بیشتر و آموزش بهتر به هر یک عملکرد مدل بهتر می‌شود.

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} = \text{input} \quad \text{filter} = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

$$\text{output size} = \text{input size} - \text{filter size} - 1$$

$$\Rightarrow 2 + 2 - 1 = 3$$

$$\text{initial output} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

برای هر ایمان فردی، فیلتر را ایمان می‌کنیم و سپس با هر سه مقیاس جمع می‌کنیم.

$$\begin{matrix} \text{ایمان (1,1)} \\ \text{filter} \times 1 \end{matrix} \rightsquigarrow \text{output} = \begin{bmatrix} 1 & 2 & 0 \\ 3 & 4 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\begin{matrix} \text{ایمان (1,2)} \\ \text{filter} \times 2 \end{matrix} \rightsquigarrow \text{output} = \begin{bmatrix} 1 & 4 & 4 \\ 3 & 10 & 8 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\begin{matrix} \text{ایمان (1,3)} \\ \text{filter} \times 3 \end{matrix} \rightsquigarrow \text{output} = \begin{bmatrix} 1 & 4 & 4 \\ 6 & 16 & 8 \\ 9 & 12 & 0 \end{bmatrix}$$

$$\begin{matrix} \text{ایمان (2,2)} \\ \text{filter} \times 4 \end{matrix} \rightsquigarrow \text{output} = \begin{bmatrix} 1 & 4 & 4 \\ 6 & 20 & 16 \\ 9 & 24 & 16 \end{bmatrix}$$



۲ = تعداد باکس هر cell  $\rightarrow$  ۷۱ yolo

$$80 \rightarrow 2 \rightarrow 90 \rightarrow 5B + C \rightarrow \text{تعداد کانال فرمبی}$$

۳ = اسکیل متفاوت  $\rightarrow$  ۳x۳ = ۹ = تعداد باکس هر cell  $\rightarrow$  ۷۳ yolo

$$255 = 3(5 + C) \rightarrow \text{تعداد کانال فرمبی}$$

۷۳ yolo  $\rightarrow$  از اسکیل های متفاوتی استفاده می کند  $\rightarrow$  object detection جعبه های بسته و اسکیل های مختلف

مختلف با نسبت جعبه های مختلف دارد، مقدار بیشتری bounding box دارد

که احتمال مکان یابی دقیق تر را افزایش می دهد.

ب) در yolo v3 از independent logistic classifier استفاده شده که هر فرایند آموزش از

bce (binary cross-entropy) به عنوان loss اضافه می کند برای داده های

که دقیقاً به یک کلاس تعلق ندارند و مورد واقع می شود. (در yolo v1 از softmax استفاده می شد که

در yolo v3 آن را غیر ضروری یافتند)

ج) از تکنیک non-maximum suppression (NMS) استفاده شده است. در این تکنیک

با توجه به ضرایب confidence و میزان Intersection over Union (IOU)

bounding box تکراری حذف می شوند و تنها جعبه با بالاترین confidence

برای یک شیء گزارش می شود.

(د) در YOLOv1، اندازه ورودی مشخص بود و scale در YOLOv2 و YOLOv3 variant

از جنبه Multi-scale training استفاده شده که در آن در فرآیند آموزش هر چند iteration

اندازه تصویر ورودی به صورت تصادفی بین  $320 \times 320$  و  $608 \times 608$  تقسیم کنند و باعث می شود

که مدل بیشتر به تغییر ابعاد تصویر ورودی حساس نباشد و به تصاویر با اندازه های مختلف به خوبی عمل کند

و scale invariant باشد.

عمدتاً این استفاده از این تکنیک امکان پذیر است این است که در YOLOv3 از bounding box

با اسکال های مختلف برای object detection استفاده می کنند.

(هـ) (۱) ابعاد مجله ها دسته های خوبی می شود. اگر چه مدل ابعاد مجله های زیادی می کند اما اگر ما

prior چیزی برای شروع به شبکه بدیم یادگیری آسان تر می شود و سریع تر هم می شود

که برای حل این مشکل ابتدا اسکالینگ و mean - k روی داده های آموزش اعمال می شود تا

ابعاد مجله و فوب به دست آورده و همچنین به جای نرم ۲ از 100 استفاده می شود.

(۲) پایداری مدل معصوم در iteration های اول (بیشتر به خاطر یافتن (۹۰۶) مجله ها بود)

راه حل ۱: از یک مکان معلوم آغاز می کنند به تدریج می آموزند که موقعیت را به روز کنند. معضات

bounding box ها نسبت به مکان cell grid مشخص می شود و از آنجمله

logistic برای قرار دادن پیش بینی ها بین ۰ و ۱ استفاده می شود.

۶ - استفاده از 53 - parknet - 53 لایه کانوئشنی و ارتباطات residual

است پس gradient بهتر دارد و ویژگی‌ها را بهتر استخراج می‌کند.  
flow

- با استفاده از independent logistic regression، می‌تواند طبقه‌بندی را برای داده‌ها

که دقیقاً یک کلاس تعلق ندارند نیز انجام دهد. (multi label classification)

- از bounding box ها با scale های مختلف استفاده می‌کنند ← در سنین

انسان با ابعاد مختلف بهتر عمل می‌کند.

- در هر scale از مقدار بیشتری bounding box استفاده می‌کنند احتمال همان یا هم دقیق‌تر

اتر است می‌باید.