

تاریخ دریافت مقاله: ۹۸/۱۲/۲۴

تاریخ پذیرش مقاله: ۹۹/۰۸/۰۵

نوع مقاله: مروری

## مروری بر شبکه عصبی پیچشی و معماری‌های مختلف آن

فاطمه باجلان

دانشکده علوم و فنون نوین - دانشگاه تهران - ایران  
پست الکترونیکی: [bajelan.fateme@ut.ac.ir](mailto:bajelan.fateme@ut.ac.ir)

هادی ویسی\*

استادیار دانشکده علوم و فنون نوین - دانشگاه تهران - ایران  
پست الکترونیکی: [h.veisi@ut.ac.ir](mailto:h.veisi@ut.ac.ir)

محمد خوانساری

دانشیار دانشکده علوم و فنون نوین - دانشگاه تهران - ایران  
پست الکترونیکی: [m.khansari@ut.ac.ir](mailto:m.khansari@ut.ac.ir)

### چکیده

مطالعه ساختارها و تکنیک‌های استفاده شده برای بهبود این شبکه می‌تواند به پژوهشگران در ایجاد ساختارهای بهینه‌تر و دقیق‌تر کمک کند.

**واژه‌های کلیدی:** یادگیری عمیق، شبکه عصبی پیچشی (کانولوشنال)، معماری‌های CNN (الکس نت، گوگل نت، رزنت)

### مقدمه

یادگیری ماشین مجموعه روش‌هایی هستند که ابداع می‌شوند تا کامپیوترها بتوانند وظایفی را به صورت خودکار یاد بگیرند. روش‌های ابتدایی یادگیری ماشین توانایی پردازش داده‌های طبیعی به فرم خام را نداشتند و نیاز به بخش مجزایی وجود داشت تا داده‌های خام تبدیل به ویژگی‌هایی شود که برای ماشین قابل استفاده باشد، برای این کار از دانش متخصصین کمک می‌گرفتند و ویژگی‌ها با الگوریتم‌هایی که مختص آن نوع داده در

شبکه عصبی پیچشی (CNN) یکی از روش‌های عمیق یادگیری ماشین است که با توجه به کارایی بالای آن، امروزه استفاده از آن در استخراج ویژگی و دسته‌بندی تصاویر امر بسیار رایجی شده است. در این مقاله برای آشنایی بیشتر پژوهشگران با این حوزه و نوآوری‌های مرتبط با این شبکه، ساختار اصلی این شبکه و تعداد ۱۰ معماری مهم و رایج آن مرور شده است. با مطالعه مقالات اصلی ارائه‌دهنده هر یک از معماری‌های شبکه پیچشی، نکات مهم و اصلی هر معماری که موجب بهبود عملکرد آن نسبت به موارد پیشین خود شده، جمع‌آوری و بررسی شده است. علاوه بر آن، مقایسه دقت معماری‌های مختلف این شبکه بخش دیگری از این مقاله مروری است. با توجه به حجم بالای مطالب در این حوزه و رشد سریع آن، نیاز به یک مرجع که همه مطالب ارائه شده تاکنون را در کنار هم مرور کرده باشد، انگیزه نوشتن این مقاله بوده است.

\* نویسنده مسئول

حوزه مشخص بود، استخراج می‌شدند تا بتواند مبنای تصمیم‌گیری ماشین باشد. اما در روش‌های جدید یادگیری ماشین، تلاش می‌شود یادگیری نحوه بازنمایی<sup>۱</sup> داده خام ورودی (تبدیل ورودی به ویژگی‌ها) نیز در خود الگوریتم گنجانده شود. در این‌گونه روش‌ها، داده می‌تواند به صورت خام به عنوان ورودی مورد استفاده قرار گیرد و نحوه بازنمایی لازم برای تصمیم‌گیری به صورت خودکار یافته می‌شود و نیازی به مهندسی ویژگی‌های آن‌ها نیست. روش‌های یادگیری عمیق<sup>۲</sup> عموماً از این دسته روش‌ها هستند که چندین لایه برای یافتن و ترکیب بازنمایی‌های مختلف از داده ورودی دارند. این لایه‌ها مستقیماً و بدون نیاز به طراحی نیروی انسانی متخصص شکل می‌گیرند و با افزایش تعداد لایه‌ها می‌توانند توابع پیچیده‌تری را یاد بگیرند. شبکه عصبی پیچشی (CNN)<sup>۳</sup> یکی از این روش‌ها است که هر دو بخش استخراج ویژگی و دسته‌بندی را در خود دارد و در صورتی که در سایر روش‌های یادگیری ماشین نیاز به ماژولی برای استخراج ویژگی از تصاویر وجود داشته باشد نیز می‌توان از بخشی از شبکه عصبی پیچشی استفاده کرد.

شبکه عصبی پیچشی عمیق کاربردهای گسترده‌ای در حوزه‌های مختلف دارد که تشخیص اشیاء در تصویر که کاربرد عام آن است و به همین جهت سالانه چالشی به نام ILSVRC<sup>۴</sup> با استفاده از دادگان ImageNet [۱] جهت تشخیص تصاویر وجود داشته که اکثر معماری‌های مهم شبکه عصبی پیچشی در روند حل این چالش ابداع شده و یا ارتقا یافته‌اند. کاربردهای خاص این شبکه‌ها در علوم مختلف نیز بسیار گسترده است و ما در اینجا به چند مورد از پژوهش‌های مرتبط به آن اشاره می‌کنیم: تشخیص چهره<sup>۵</sup> [۲]، پردازش متن [۳]، در علم پزشکی تشخیص بیماری‌های مختلف [۴] مانند سرطان سینه [۵]، بیماری‌های دندانی [۶]، و سایر بیماری‌هایی که به کمک تصاویر قابل تشخیص است. در علم کشاورزی

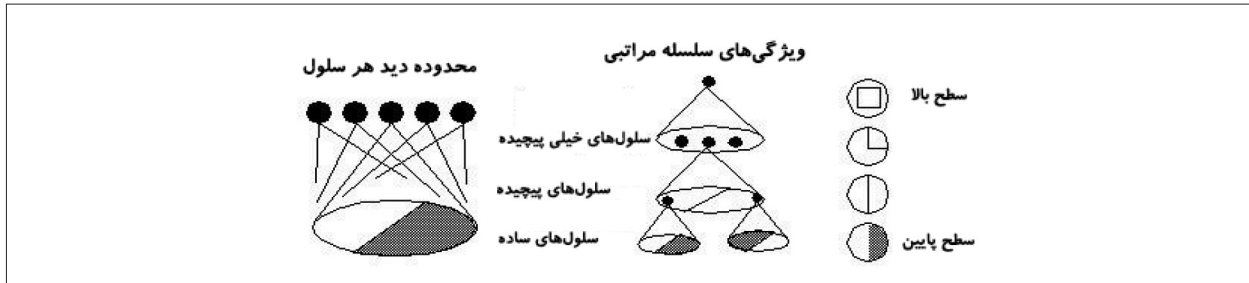
1- Representation  
2- Deep Learning  
3- Convolutional Neural Networks (CNN)  
4- ImageNet Large Scale Visual Recognition Challenge (ILSVRC)  
5- Face Recognition

تشخیص بیماری‌ها یا کمبودهای شایع گیاهان در سطح کلان و با کمک تصاویر ماهواره‌ای [۷] و یا در سطح خرد به کمک تصاویر معمول [۸]. در علوم محیط زیست تشخیص الگوهای بوم‌شناسی [۹] و حتی در حوزه‌های مرتبط به بازی و بازیکن هوشمند مانند انجام بازی گو<sup>۱۰</sup> [۱۰] از شبکه‌های عصبی پیچشی استفاده شده است.

افزایش حجم دادگان و گستردگی نیاز به پردازش و دسته‌بندی تصاویر، نیاز به افزایش دقت در این حوزه را ایجاد کرده و هر ساله در دنیا معماری‌های مختلفی برای شبکه عصبی پیچشی ارائه می‌شود. با توجه به کاربردهای گسترده این‌گونه شبکه‌ها، نیاز به اطلاع از آن‌ها و مرور تکنیک‌های استفاده شده برای بهبود کارکردشان مورد نیاز است و انگیزه این مقاله ارائه مرجعی برای معرفی علمی این شبکه و انواع توسعه‌های مرتبط با آن در کنار هم بوده است. در ادامه مقاله، علاوه بر معرفی شبکه عصبی پیچشی و لایه‌های رایج در آن در بخش دوم، معماری‌های مهم و ویژگی‌های اصلی آن‌ها در بخش سوم مرور شده است. در بخش چهارم تکنیک‌های مورد استفاده برای بهبود معماری‌ها شرح داده شده و در فصل آخر نیز بر اساس پژوهش‌های موجود، مقایسه‌ای میان معماری‌های ذکر شده آمده است.

## ۲- ساختار شبکه عصبی پیچشی

موضوع شبکه‌های عصبی پیچشی از سال ۱۹۹۰ مطرح بوده است [۱۱]، اما دلایل اصلی این‌که شبکه‌های عصبی پیچشی در چند سال اخیر امکان استفاده و ارتقاء یافته‌اند و به روشی غالب در فضای یادگیری ماشین تبدیل شده‌اند، افزایش حجم داده در دسترس برای آموزش و همچنین توسعه توان پردازشی کامپیوترها بوده است. این دو موضوع امکان افزایش تعداد لایه‌های شبکه‌های عصبی را ایجاد کرده تا بتوان ویژگی‌های پیچیده‌تری را از داده‌ها استخراج نموده و بر آن اساس دسته‌بندی را با دقت بالاتر انجام داد.



شکل ۱: ساختار سلسله مراتبی سیستم بینایی و مغز انسان برای تشخیص اشیاء [۱۳]

تفاوت این شبکه با سایر شبکه‌های عصبی این است که ورودی آن یک تصویر در نظر گرفته شده و روش انجام اعمال بر اساس ورودی تصویر تطبیق پیدا کرده‌اند و ساختاری را پیشنهاد می‌دهد که بتوان بهترین ویژگی‌ها را از تصاویر استخراج کرد. در ساختار این شبکه ورودی سه‌بعدی در نظر گرفته می‌شود؛ تصویری که دارای عرض، ارتفاع، و عمق است (شکل ۲).

یان لکان<sup>۱۲</sup>، به‌عنوان یکی از اولین افرادی که در خصوص شبکه‌های عصبی پیچشی صحبت کرده است [۱۵]، چند دلیل را برای رسیدن به چنین شبکه‌ای عنوان می‌کند: ۱. کاهش تعداد پارامترها برای کاهش پیچیدگی شبکه. ۲. استفاده از ویژگی محلی بودن مقادیر در تصاویر. ۳. عدم وابستگی ویژگی‌ها به یک نقطه خاص از تصویر، که تحت عنوان اشتراک پارامترها مطرح می‌شود. در ادامه ساختار این شبکه و لایه‌های مختلف آن مرور می‌شوند.

## ۲-۱- لایه‌های شبکه عصبی پیچشی

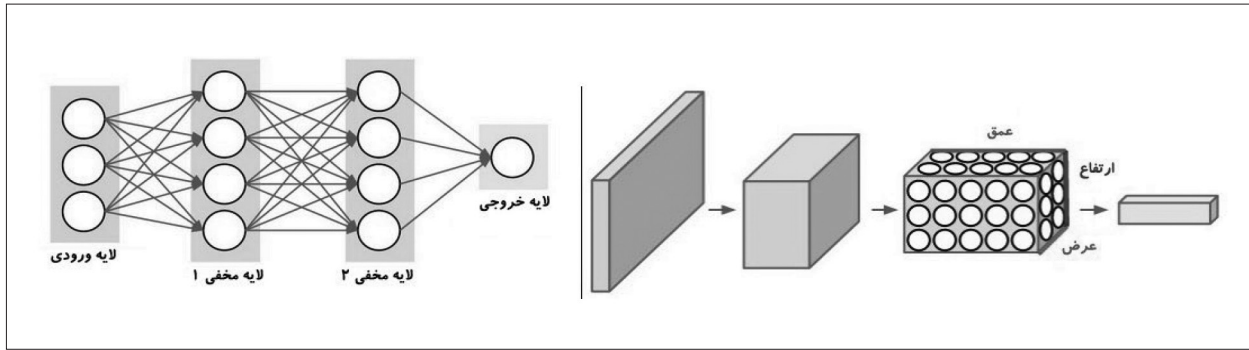
شبکه عصبی پیچشی مانند سایر شبکه‌های جلودر از تعدادی لایه متوالی تشکیل شده است که علاوه بر ورودی و خروجی، لایه‌های اصلی برای تشکیل این شبکه لایه پیچش<sup>۱۳</sup>، لایه ادغام<sup>۱۴</sup> (یا نمونه‌برداری)، لایه تمام متصل (مشابه آنچه در شبکه‌های عصبی پرسپترون چند لایه<sup>۱۵</sup> داریم) هستند که همراه آن‌ها موارد دیگری هستند که گاهی لایه خوانده می‌شوند، در ادامه مورد بررسی قرار می‌گیرند.

شبکه عصبی پیچشی، نوعی از شبکه‌های عصبی مصنوعی جلودر<sup>۷</sup> است که برای پردازش و دسته‌بندی داده‌های چندبعدی مانند تصاویر ایجاد شده است. ساختار آن در تلاش برای شبیه‌سازی کارکرد سیستم بینایی انسان به‌صورت فعلی در آمده است. همان‌طور که در شکل ۱ نشان داده شده است، بر اساس نتایج پژوهش‌های هیوبل و ویزل<sup>۸</sup> در سال ۱۹۵۹ سلول‌های چشم و مغز برای تشخیص اشیاء یک ساختار سلسله‌مراتبی از ویژگی‌ها را از ساده تا پیچیده تشخیص می‌دهند [۱۲]، بدین گونه که ساده‌ترین ویژگی‌های تصاویر مانند مرز به‌عنوان ویژگی‌های لایه اول و در لایه‌های بالاتر مواردی مانند انحنا و زاویه و در ادامه شکل‌ها و ویژگی‌های پیچیده‌تر تشخیص داده می‌شوند. علاوه بر آن، ویژگی دیگر این شبکه‌ها این است که هر سلول، نه تنها یک نقطه بلکه یک محدوده دوطرفه را مشاهده می‌کند، که به آن ناحیه دریافت<sup>۹</sup> می‌گویند. این نواحی دریافت که اطلاعات همسایگی هر نقطه را تامین می‌کنند، غالباً با هم هم‌پوشانی دارند.

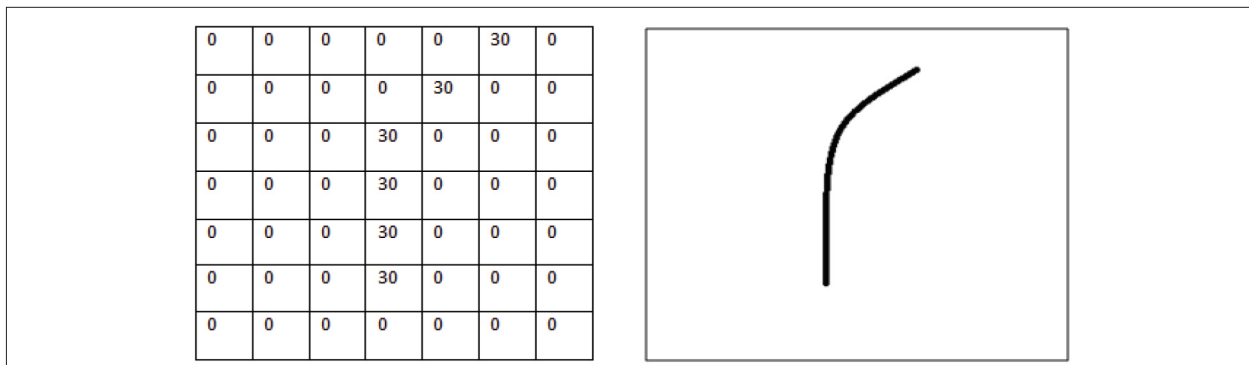
شبکه‌های عصبی پیچشی مانند سایر شبکه‌های عصبی مصنوعی، از نورون‌هایی ساخته شده‌اند که وزن‌ها و بایاس<sup>۱۰</sup>‌ها را ذخیره می‌کنند و در لایه‌های بعد خود روشی برای تصمیم‌گیری را به همراه دارند. اگر بخواهیم از شبکه‌های عصبی معمولی و لایه‌های تمام متصل<sup>۱۱</sup> برای دسته‌بندی تصاویر استفاده کنیم حجم هر لایه مخفی بسیار بزرگ و روند به‌روزرسانی آن بسیار طولانی خواهد بود.

7- Feedforward  
8- Hubel and Wiesel  
9- Receptive Field  
10- Bias  
11- Fully Connected

12- Yann Lecun  
13- Convolution  
14- Pooling  
15- Multi-Layer Perceptron (MLP)



شکل ۲: ساختار کلی شبکه‌های عصبی (چپ) - ساختار تغییر یافته برای ورودی‌های تصویری یا چندبعدی (راست) [۱۴]



شکل ۳: بازنمایی یک پالایه تشخیص دهنده انحنای [۱۶]

## ورودی

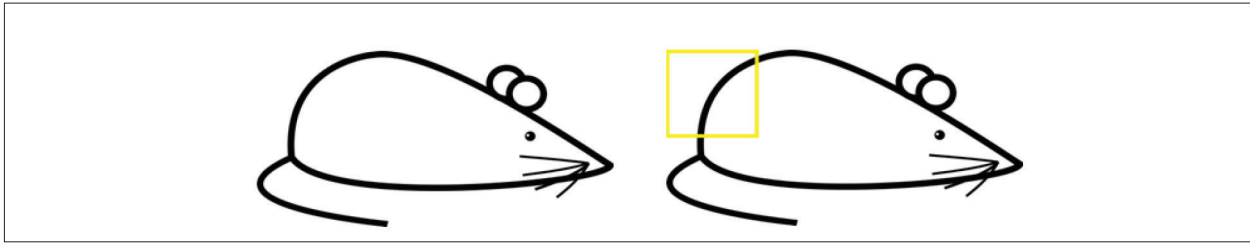
ورودی این شبکه یک تصویر است که یک داده سه‌بعدی دارای طول و عرض و عمق به شمار می‌رود. عمق به معنای تعداد مؤلفه‌هایی است که با آن مقدار هر پیکسل را نمایش می‌دهیم، مثلاً در صورتی که تصویر سیاه و سفید باشد عمق برابر ۱ و در صورتی که تصویر رنگی با سه مؤلفه‌ی قرمز، سبز و آبی باشد، عمق برابر ۳ خواهد بود و اگر در فضای دیگری رنگ‌ها را نمایش دهیم عمق برابر تعداد مؤلفه‌های آن فضا خواهد بود. مقادیر پیکسل‌ها در داده‌های تصویری طبیعی به صورت محلی با هم ارتباط دارند، یعنی اگر ناحیه‌ای دارای رنگ خاصی باشد، پیکسل‌های اطراف آن نیز دارای رنگ مشابه خواهند بود. به همین دلیل حفظ ساختار مقادیر در تصویر و پردازش آن‌ها در ساختار اصلی حائز اهمیت است. در عمل، بخشی از تصویر ورودی معادل ورودی شبکه است که یک پنجره محدود از تصویر است که روی تصویر لغزانده می‌شود تا کل تصویر را پیمایش کند.

## پیچش

این لایه شبیه‌سازی عملکرد سلول‌های چشم انسان را بر عهده دارد و در واقع لایه‌ای است که ویژگی‌های شبکه را تشخیص می‌دهد و در خود ذخیره می‌کند. در این لایه مفهومی به نام پالایش تعریف می‌شود که مجموعه نورون‌های لایه مخفی شبکه عصبی هستند، با این تفاوت که در شبکه عصبی مصنوعی هر نورون مخفی متصل به کل ورودی‌های لایه قبل خود یا بخش مشخصی از آن‌هاست، اما در لایه پیچش، این پالایش‌ها به بخش خاصی اختصاص ندارند و یک ویژگی را در کل ماتریس جست‌وجو می‌کنند (مانند آنچه در شکل‌های ۳ و ۴ نشان داده شده است).

یکی از مزایای محاسباتی این ساختار، اشتراک پارامترها<sup>۱۶</sup> است که به این دلیل رخ می‌دهد که پالایه‌ها به ازای کل نواحی تصویر ساخته و به‌روزرسانی می‌شوند. مزیت دیگر این ساختار به‌عنوان اتصال تنگ<sup>۱۷</sup> شناخته

16- Parameter Sharing  
17- Sparse Connectivity



شکل ۴: بازنمایی از تطبیق پالایه بر انحنای تصویر [۱۶]

مؤثر هستند، عمق یا همان تعداد پالایه‌ها، گام و حاشیه، که در ادامه شرح هر یک آمده است.

**عمق:** این پارامتر نشان‌دهندهٔ تعداد پالایه‌هاست. هر پالایه به یک ویژگی در تصویر اختصاص دارد، در نتیجه افزایش تعداد پالایه‌ها ویژگی‌های استخراج‌شده را افزایش می‌دهد. پالایه‌ها اختصاص به یک لایه از تصویر ندارند و از لحاظ عمق، هم عمق ورودی خود هستند، همانطور که چشم انسان نیز لایه‌های رنگی را مجزا از یکدیگر نمی‌بیند. **گام:**<sup>۲۰</sup> مشخص‌کننده اندازه گام در پیچش پالایه بر روی تصویر است و در واقع بیانگر تعداد پیکسل‌هایی است که در لغزاندن پنجره ورودی روی تصویر در هر مرحله، پنجره را جابجا می‌کنیم. در حالت پیش‌فرض این عدد ۱ است، یعنی پالایه یک پیکسل یک پیکسل حرکت می‌کند. هر چه مقدار گام بزرگ‌تر باشد اندازه خروجی کوچک‌تر خواهد بود. گام هم در راستای محور طول هست و هم در راستای محور عرض. پنجره ورودی ابتدا در راستای افقی روی تصویر لغزانده می‌شود و زمانی که به انتهای آن ردیف از تصویر رسید به ابتدای ردیف بعدی با اندازه گام عمودی جابجا می‌شود.

**حاشیه:**<sup>۲۱</sup> یکی دیگر از پارامترهای مؤثر در اندازه خروجی، عددی است که با عنوان حاشیه تعریف می‌شود. این متغیر مشخص‌کننده تعداد ردیف صفرهایی است که به حاشیه تصویر ورودی اضافه می‌کنیم. این پارامتر برای تعیین اندازه تصویر خروجی از لایه پیچش مهم است. با فرض این‌که اندازه تصویر ورودی  $W$ ، اندازه پالایه  $F$ ، اندازه گام  $S$  و اندازه حاشیه  $P$  باشد، آنگاه اندازه تصویر

شده است. کاهش تعداد اتصال‌ها - نسبت به اتصال کامل - و ایجاد نظم بخصوص، موجب کاهش تعداد پارامترها و در نتیجه کاهش محاسبات می‌شود.

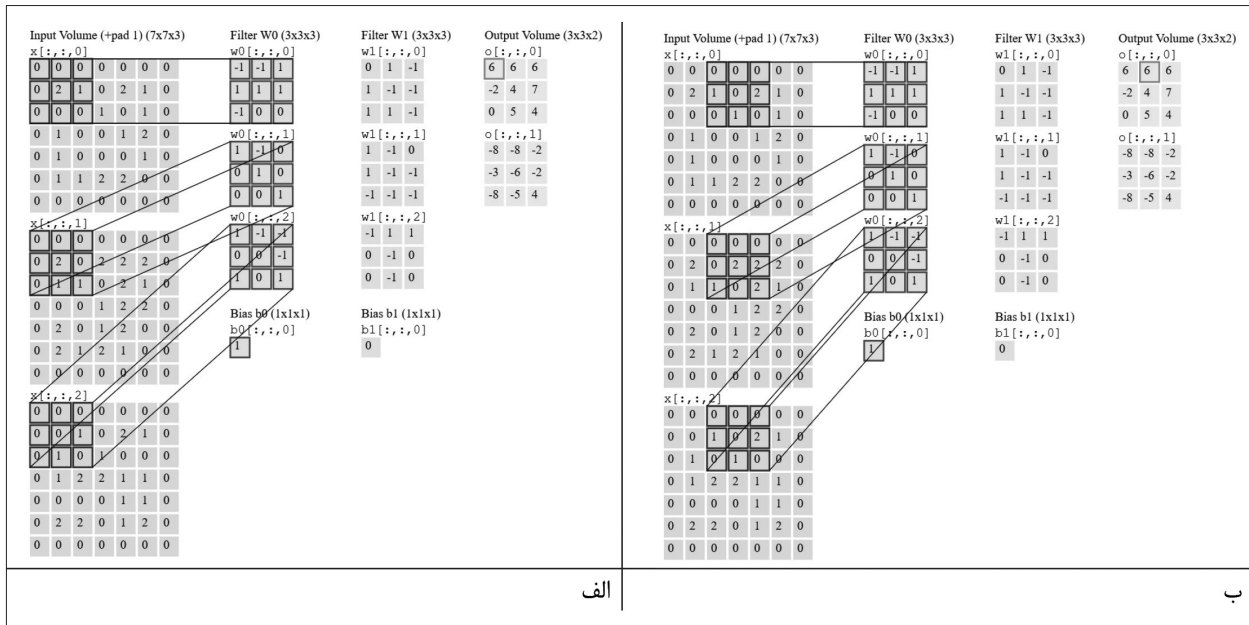
پالایه‌ها، یا همان پارامترهای لایه پیچش، ماتریس‌هایی سه‌بعدی هستند که وزن‌های شبکه در این لایه را در خود ذخیره می‌کنند. این وزن‌ها باید در مرحله آموزش از روی داده آموزشی تخمین زده شوند. این پالایه‌ها با پیچش در ورودی لایه خود، در تمامی ابعاد، یک خروجی دو بعدی تولید می‌کنند که آن را نگاشت فعال‌سازی<sup>۱۸</sup> یا نگاشت ویژگی<sup>۱۹</sup> می‌نامیم. در واقع هر پالایه، پنجره ورودی از تصویر را پالایش کرده و آن پنجره از تصویر به یک عدد تبدیل می‌گردد. اندازه هر پالایه که برابر با اندازه پنجره لغزان روی تصویر است، از جمله پارامترهای شبکه است که توسط طراح شبکه تعیین می‌شود. بعد از یادگیری، هرکدام از این پالایه‌های لایه پیچش بیانگر یک ویژگی از تصویر می‌شوند. به همین دلیل ما در هر لایه بیش از یک پالایه داریم، و عمق خروجی یک لایه همان تعداد پالایه‌هاست. هر پنجره از ورودی با عبور از هر پالایه (ضرب شدن مقادیر متناظر در همدیگر و جمع شدن) یک پیکسل از تصویر خروجی ایجاد می‌کند که این موضوع در رابطه (۱) نشان داده شده است. در این رابطه  $I_j^l$  بیانگر ویژگی (پیکسل)  $j$ م از تصویر لایه  $l$ ،  $w_{ij}^l$  وزن‌های پالایه لایه  $l$ ، و  $b_j^l$  بایاس است. در اینجا  $f(\cdot)$  تابع فعال‌سازی است که معمولا ReLU است.

$$I_j^l = f\left(\sum_i I_i^{l-1} \times w_{ij}^l + b_j^l\right) \quad (1)$$

در خروجی لایه پیچش، بجز اندازه پالایه سه عامل

20- Stride  
21- Padding

18- Activation Map  
19- Activation Map



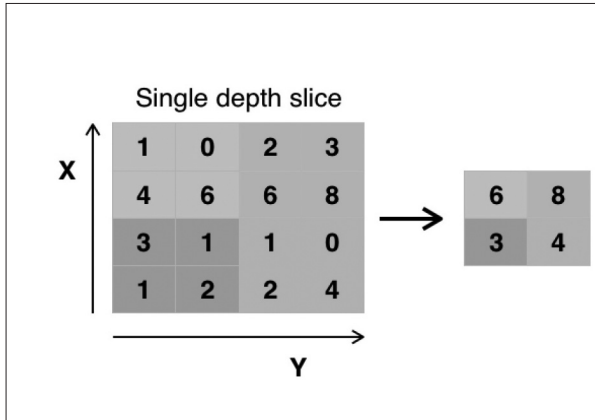
شکل ۵: مثال پیچش پالایه ۳ در ۳ در یک ورودی ۵ در ۵ با اندازه گام ۲ و مقدار حاشیه ۱ [۱۴]

به کمک توابع فعال‌سازی مختلف ایجاد می‌شود زیرا اکثر دسته‌بندی‌های دنیای طبیعی دارای مرزهای خطی مشخص نیستند و ما نیاز به شبیه‌سازی توابع به صورت غیرخطی داریم. در واقع ReLU یک تابع فعال‌سازی است که روی مقادیر پیکسل‌های تصویر حاصل از خروجی لایه پیچش اعمال می‌شود و می‌توانست به جای آن هر تابع فعال‌سازی غیرخطی دیگر مانند سیگموید هم باشد اما یکی از مزایای استفاده از ReLU نسبت به توابعی مانند سیگموید این است که محاسبه آن بسیار سریع‌تر است، و این موجب می‌شود استفاده از آن برای آموزش شبکه‌های عمیق کم‌هزینه‌تر و رایج‌تر باشد. نکته دیگر این است که آموزش شبکه با سیگموید به سرعت به حد اشباع می‌رسد، به این دلیل که گرادینت‌ها تنها در ناحیه کوچکی، غیرصفر و تأثیرگذار در به‌روزرسانی هستند، و در سایر نقاط مقادیری بسیار نزدیک به صفر می‌شوند که ضرب شدن آن‌ها در هم موجب می‌شود عملاً به گرادینت صفر برسیم، و دیگر شاهد تغییراتی در نورون‌ها نباشیم. به این دلایل، ReLU به‌عنوان تابع فعال‌سازی رایجی در این شبکه‌ها و برخی شبکه‌های عمیق دیگر تبدیل شده است. رابطه و نمودار این تابع و مقایسه آن با دو تابع فعال‌سازی رایج

خروجی  $(W-F+2P)/S+1$  است. این متغیرها می‌توانند اندازه طول یا عرض تصویر باشند. بنابراین، به صورت عمومی اثر پالایش کردن لایه پیچش، تولید ماتریس‌هایی با اندازه کوچک‌تر از تصویر ورودی است و برای حفظ اندازه و یا تولید اعداد صحیح (غیراعشاری) برای اندازه ماتریس خروجی این لایه از پارامتر حاشیه استفاده می‌شود. در شکل ۵ یک مثال آورده شده است که در آن اندازه تصویر ورودی  $7 \times 7$  و رنگی (با عمق ۳) است و اندازه هر پالایه (پنجره روی تصویر)  $3 \times 3$  است؛ دو بار عمل پیچش در یک پالایه نشان داده شده است. اندازه گام در این مثال ۲ و مقدار حاشیه ۱ است. در این مثال، اندازه تصویر خروجی  $3 \times 3$  است و به تعداد پالایه‌ها (که در این مثال دو پالایه است) ماتریس خروجی داریم.

### ReLU

اگر لایه را، برای مفهومی استفاده کنیم که در هر بار آموزش دارای فرآیند به‌روزرسانی باشد، ReLU رانمی‌توان یک لایه به حساب آورد و صرفاً می‌توان به دید یک تبدیل بر روی مقادیر خروجی لایه پیچش به آن نگاه کرد، تبدیلی که ویژگی خطی بودن را در شبکه از بین می‌برد. از بین رفتن این ویژگی مزیتی است که در شبکه‌های عصبی



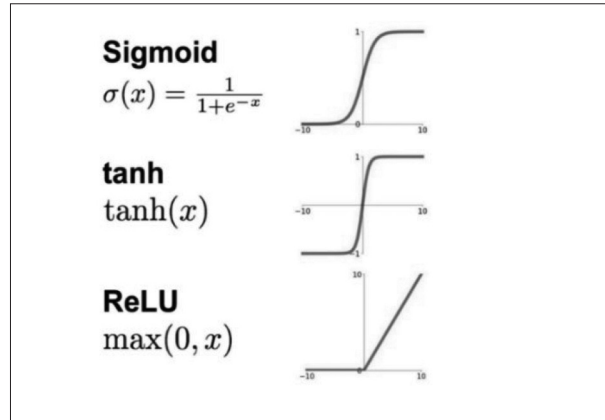
شکل ۷: عملکرد لایه ادغام با عملگر بیشینه‌گیری با پالایه ۲×۲ و گام ۲

### ۲۳ نرمال‌سازی دسته‌های

نرمال‌سازی دسته‌های مقادیر ماتریس‌ها را که بازه نامشخصی دارند، به بازه مشخصی مانند ۰ تا ۱ تبدیل می‌کند. برای این کار مزایای متفاوتی ذکر شده است: افزایش سرعت یادگیری شبکه، افزایش قدرت و کاهش بیش‌برازش در شبکه. البته با توجه به محل قرارگیری آن در شبکه و لایه‌هایی که پیش و پس از آن قرار می‌گیرند برخی از این مزایا را ایجاد می‌کند. این عملگر هم که نمی‌توان آن را لایه نامید، فقط در برخی از معماری‌ها استفاده می‌شود.

### تمام متصل

لایه تمام متصل که شاکله اصلی شبکه‌های عصبی معمولی مانند پرسپترون چندلایه است، عموماً در انتهای ساختار شبکه عصبی پیچشی به‌عنوان دسته‌بند به کار گرفته می‌شود. این لایه خود شامل یک لایه خروجی و یک یا چند لایه ماقبل آن (مشابه لایه مخفی در شبکه پرسپترون چندلایه) است. هرکدام از این لایه‌ها می‌توانند تابع فعال‌سازی متفاوتی مانند تانژانت هایپربولیک و سیگموید داشته باشند. تابع فعال‌سازی لایه آخر معمولاً بیشینه‌نرم (Softmax) است که همان‌طور که در رابطه (۲) آمده است، این لایه مشابه سیگموید است با این تفاوت که مقادیر خروجی آن به‌صورت احتمال است و برای کاربردهای دسته‌بندی که دنبال یافتن دسته‌ای با بیشترین احتمال برای ورودی شبکه هستیم، کاربرد فراوانی دارد. با توجه به استفاده زیاد شبکه



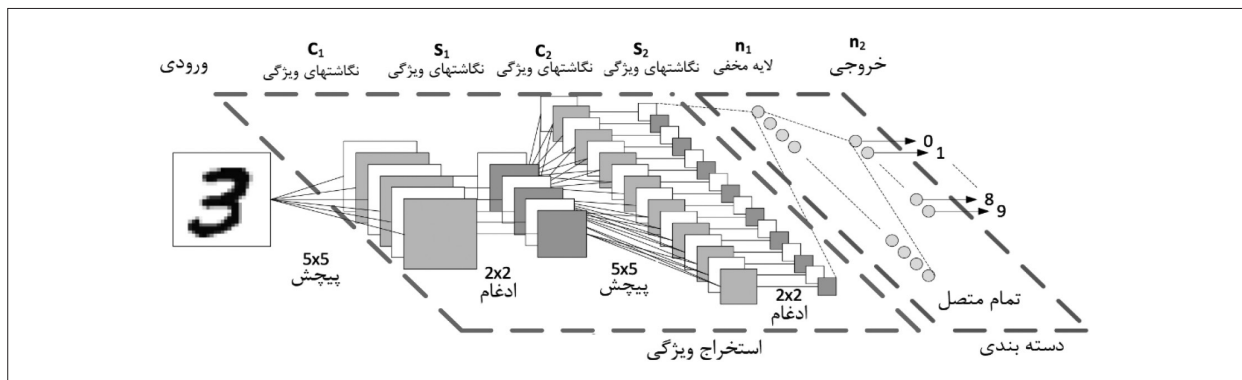
شکل ۶: تابع فعال‌سازی ReLU و دو تابع معمول دیگر در شبکه عصبی عمیق

دیگر (سیگموید و تانژانت هایپربولیک) در شکل ۶ آورده شده است.

محاسبات لایه‌های بعد از اعمال این تابع فعال‌سازی، که می‌تواند مربوط به یک لایه ادغام باشد و یا یک لایه پیچش دیگر، بر روی ماتریس‌های خروجی عبور داده شده از این تابع اعمال می‌شود.

### ادغام

کاربرد این لایه، کاهش اندازه نقشه‌های ویژگی، بدون افزودن پارامتری به شبکه است. برای کاهش تعداد پارامترها و کاهش میزان پردازش در شبکه از لایه ادغام استفاده می‌شود. ادغام به‌صورت مجزا بر هر لایه از نگاشت فعال‌سازی اعمال می‌شود. هدف مهم دیگری که از این کار دنبال می‌شود جلوگیری از بیش‌برازش ۲۲ شبکه است. این لایه معادل با اعمال یک پالایه (معمولاً با اندازه ۲×۲) روی ماتریس حاصل از لایه پیچش برای کاهش اندازه است و عملگر پالایه بیشینه‌گیری یا میانگین‌گیری است. این پالایه با گام (معمولاً برابر با ۲) جلو می‌رود و با این کار طول و عرض ماتریس کوچک می‌شود. بدیهی است در این لایه پارامتری برای یادگیری (مانند وزن‌های شبکه) وجود ندارد. در شکل ۷ یک مثال از عملکرد این پالایه ۲×۲ و با گام ۲ نشان داده شده است که باعث شده اندازه تصویر نصف شود.



شکل ۸: نمونه‌ای از یک شبکه عصبی پیچشی برای دسته‌بندی اعداد دست نوشته

در زیربخش استخراج ویژگی در این شبکه، چیدمان دو لایه پیچش و ادغام و تعداد آن‌ها از جمله پارامترهایی است که باید در زمان طراحی شبکه در مورد آن‌ها تصمیم‌گیری شود. ساختار کلی چیدمان این لایه را می‌توان به صورت زیر بنویسیم که بیان می‌کند N لایه پیچش متوالی که بعد از هر کدام از آن‌ها تابع ReLU اعمال شده، یک لایه ادغام به صورت اختیاری (بدین معنی که می‌تواند وجود نداشته باشد) می‌آید و کل این ساختار می‌تواند M بار تکرار شود.

$$[[\text{CONV} \rightarrow \text{RELU}] * N \rightarrow \text{POOL?}] * M$$

در زیربخش دسته‌بندی نیز، لایه خروجی یک بار و لایه (مخفی) ماقبل آن می‌تواند به تعداد K بار (حداقل یک بار) با تابع فعال‌سازی غیرخطی مرتبط (مانند تابع ReLU) تکرار شود. لذا کل ساختار این شبکه را می‌توان به صورت زیر نوشت.

$$\text{INPUT} \rightarrow [[\text{CONV} \rightarrow \text{RELU}] * N \rightarrow \text{POOL?}] * M \rightarrow [\text{FC} \rightarrow \text{RELU}] * K \rightarrow \text{FC}$$

برای آموزش شبکه‌های عصبی پیچشی که معادل با تنظیم و به‌روزرسانی مقادیر پالایه‌ها در لایه‌های پیچش و وزن‌های لایه‌های تمام متصل است، مشابه یادگیری شبکه‌های عصبی پرسپترون چندلایه از الگوریتم یادگیری پس انتشار خطا<sup>۲۴</sup> استفاده می‌شود [۱۱] [۲۰]. در این الگوریتم، ابتدا نمونه داده آموزش به شبکه داده می‌شود تا خروجی آن محاسبه شود، سپس با استفاده از یک تابع خطا مانند میانگین مربعات خطا<sup>۲۵</sup> یا آنتروپی متقابل<sup>۲۶</sup> خطای شبکه

پیچشی در کاربردهای دسته‌بندی، از این تابع در لایه آخر این شبکه استفاده می‌شود.

$$s(x_i) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} \quad (2)$$

در نهایت ساختار یک شبکه عصبی پیچشی به صورت کلی شامل دو زیربخش استخراج ویژگی و دسته‌بندی است که در زیربخش استخراج ویژگی از تعدادی لایه پیچش و ادغام استفاده می‌شود و زیربخش دسته‌بندی که در واقع همان پرسپترون چند لایه است، فقط شامل لایه‌های تمام متصل است. شکل ۸ مثالی از ساختار کلی این شبکه برای دسته‌بندی تصاویر اعداد دست‌نوشته (۱۰ دسته) با دو لایه پیچش و ادغام متوالی و لایه تمام متصل با یک لایه خروجی و یک لایه مخفی به همراه اندازه‌های تصاویر هر لایه آورده شده است. همان‌گونه که پیداست، اندازه تصاویر در هر لایه کوچک شده اما تعداد آن‌ها بیشتر می‌شود. در واقع هر لایه مسئول استخراج سطحی از انتزاع ویژگی‌های ماتریس ورودی را برعهده دارد و هرچه از لایه ورودی دورتر می‌شویم و به سمت لایه خروجی حرکت می‌کنیم، سطح انتزاع بیشتر می‌شود؛ این مفهوم کارکرد اصلی استخراج ویژگی در شبکه عصبی عمیق است. در آخرین لایه زیربخش استخراج ویژگی (که در شکل ۸ یک لایه ادغام است)، همه مقادیر در کنار هم تبدیل به یک بردار شده و به‌عنوان ویژگی استخراج شده توسط شبکه پیچشی به لایه تمام متصل داده می‌شوند تا روی آن‌ها دسته‌بندی صورت گیرد.

24- Error Back-Propagation  
25- Mean Square Error (MSE)  
26- Cross Entropy



	1	2	3	4	5	6	7	8	9	10	11	12
1	X	X	X		X	X						
2		X	X	X								
3				X	X		X	X	X	X		
4										X	X	X

شکل ۹: اتصالات بین لایه H2 و H3 در شبکه لینت [۱۱]

را متناسب با این ویژگی نظم داده است. روش‌های کلاسیک تشخیص الگو در تصاویر، عموماً ویژگی‌های سطح پایین را تشخیص داده و با ترکیب آن‌ها باهم به ویژگی‌های سطح بالاتر دست پیدا می‌کنند. در این ساختار نیز با همین ایده، دو لایه پیچش قرار داده شده است که برای دستیابی به ویژگی‌های گفته‌شده، نورون‌های این لایه‌ها، اطلاعات یک همسایگی از ماتریس ورودی را در خود دارند.

این شبکه دارای چهار لایه، H1 و H2 و H3 و H4 است. که لایه اول و سوم لایه‌های پیچش با ماتریس وزن‌های ۵×۵ هستند، و لایه‌های دوم و چهارم لایه‌های نمونه‌برداری<sup>۲۹</sup> هستند - که بعداً به نام لایه ادغام شناخته شدند - و میانگین‌گیری آن‌ها میانگین وزنی، روی نواحی ۲×۲ است. ورودی لایه H1، تصویر اصلی است، از ۴ پالایه مورد استفاده، ۴ نقشه ویژگی<sup>۳۰</sup> - البته در آن زمان واژه پالایه برای این کاربرد استفاده نمی‌شده است - به دست می‌آید. ورودی H3، ۴ ماتریس است که نقشه‌های ویژگی به دست آمده از دو لایه قبلی هستند و با توجه به محدودیت توان پردازشی کامپیوترها در آن سال‌ها، این لایه ۱۲ پالایه یک یا دو لایه دارد و این تعداد برای پوشش دادن تمام ترکیبات نقشه‌های ویژگی کم است. به همین جهت هر کدام از پالایه‌ها روی یک یا دو تا از نقشه ویژگی‌ها اعمال می‌شوند. شکل ۷ تطبیق پالایه‌ها به نقشه‌های ویژگی لایه H2 را نشان می‌دهد.

این شبکه که برای تشخیص اعداد دست‌نوشته طراحی شده بود، در انتهای شبکه، مستقیماً لایه تصمیم‌گیری را داریم که دارای ۱۰ نورون است (برابر با تعداد ۱۰ رقم) و تمامی واحدهای خروجی لایه H4 به آن اتصال کامل دارند. تابع فعال‌سازی که پس از ادغام اعمال می‌شود، تانژانت هایپربولیک، و تابع فعال‌سازی لایه آخر، سیگموید است.

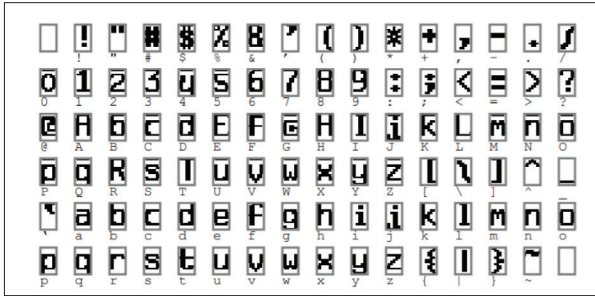
(در مقایسه با خروجی واقعی) محاسبه شده و با به‌روز کردن پارامترهای شبکه در جهت کاهش گرادیان<sup>۳۷</sup> سعی در کمینه کردن خطا می‌شود. برای انجام این کار، در هر لایه لازم است خطا محاسبه شود که به این منظور خطای لایه آخر به لایه ماقبل خود پس انتشار داده می‌شود تا خطای لایه ماقبل آخر به دست آید و از پس انتشار خطای این لایه، خطای مربوط به لایه قبل‌تر و این فرایند تا رسیدن به لایه ورودی ادامه پیدا می‌کند.

### ۳- معماری‌های شبکه عصبی پیچشی

بعد از مورد توجه قرار گرفتن شبکه عصبی پیچشی از سال ۲۰۱۲ به بعد، توسعه‌های مختلفی روی آن‌ها انجام شده است که منجر به معماری‌های نوینی از این شبکه‌ها شده است. این معماری‌ها هرچند در مفاهیم پایه بیان شده در بخش پیشین مشترکند اما هر کدام دارای ایده‌هایی برای بهبود کارایی هستند. در این بخش تعداد ۱۰ مورد از معماری‌های مهم شبکه عصبی پیچشی مرور می‌شوند و ایده‌های پر استفاده‌تر و تکرار شده‌ی آن‌ها در فصل بعد به تفکیک قابل مطالعه است.

#### ۱- لینت<sup>۲۸</sup>

یان لکان که ارائه‌دهنده اولین معماری از شبکه عصبی پیچشی است، در ابتدا شبکه لینت-۱ [۱۱] را در سال ۱۹۹۰ معرفی نمود. فرض مطرح شده برای طراحی این ساختار این است که یک شبکه عصبی تمام متصل با ورودی تصویر، به دلیل تعداد بالای پارامترها، قدرت تعمیم پایینی دارد، پس نیاز است شبکه‌ای با اتصالات محدودتر طراحی شود. همچنین با این فرض که ورودی‌ها تصویر هستند، از ویژگی محلی بودن داده‌ها استفاده شده و کاهش پارامترها



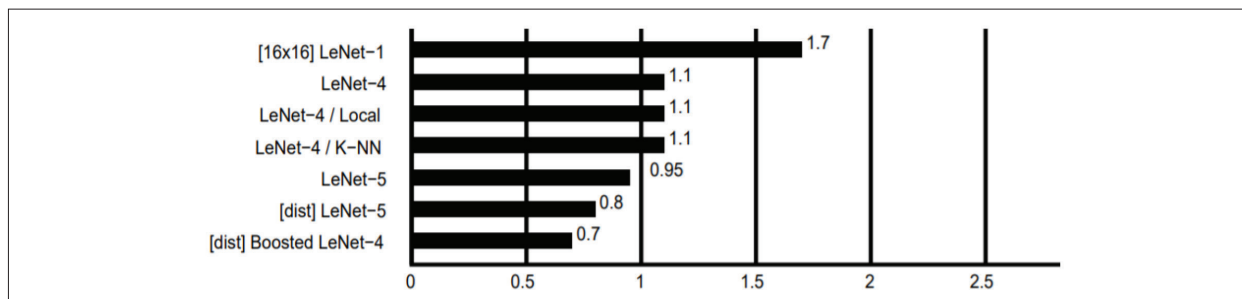
شکل ۱۰: نویسه‌های ساده شده با اندازه ۷×۱۲ به عنوان مرجع محاسبه فاصله برای خروجی لینت ۵ [۱۷]

شبکه است، که با اعمال تبدیل‌هایی مانند چرخش، مقیاس کردن و جابجایی، تعداد نمونه‌های مجموعه دادگان افزایش می‌یابد.

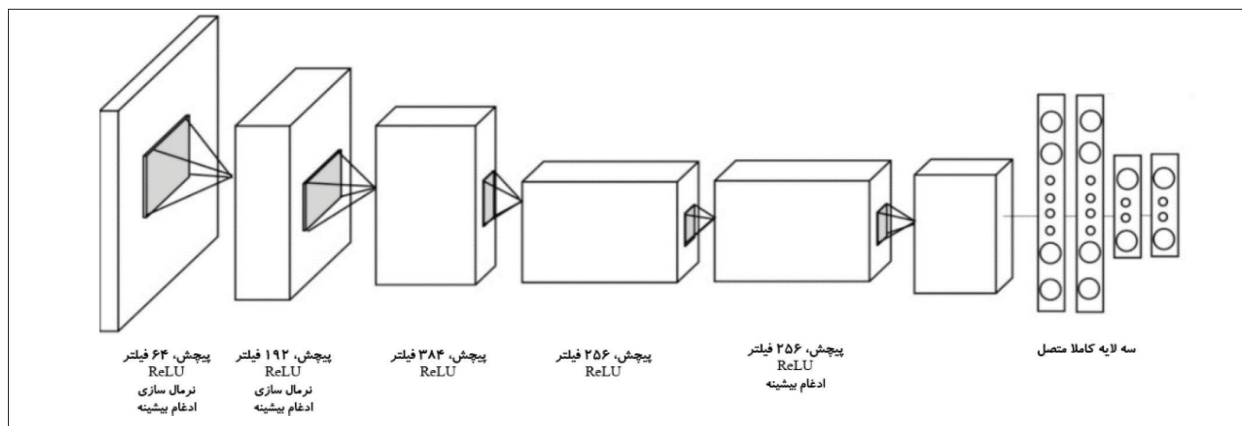
لینت-۵، [۱۷] که نسخه ارتقایافته لینت مربوط به سال ۱۹۹۵ است و در مقالات متعددی به آن اشاره شده و مرجع مقایسه قرار گرفته است، از لحاظ ساختاری شباهت زیادی به لینت-۴ دارد. موارد تفاوت آن در تعداد لایه‌های پنهان پیش از لایه تصمیم‌گیری است که یک لایه افزایش داشته است، یعنی در انتهای شبکه، دو لایه مخفی تمام متصل ۱۲۰ تایی و ۸۴ تایی قرار دارند. همچنین در لایه اول پیچش، ۶ نقشه ویژگی دارد که این تعداد در لینت-۴، برابر با ۴ نقشه بود. این تغییرات کمک می‌کند میزان ویژگی‌های استخراج‌شده در لایه اول افزایش یابد و از ترکیب آن‌ها در لایه سوم ویژگی‌های متنوع‌تری به دست آید. در این شبکه برای حرکت از لایه دوم - که در آن ادغام و اعمال تابع فعال‌سازی می‌شود - به لایه سوم، از ۶ نقشه ویژگی به ۱۶ نقشه ویژگی می‌رسیم که برای این کار، طبق روال مشخصی سه، چهار یا شش نقشه ویژگی از لایه دوم با هم‌دیگر ترکیب شده‌اند. در این ساختار، عمل ادغام به این شکل انجام می‌شود که ابتدا مقادیر داخل ناحیه ۲×۲ با هم جمع شده و سپس در یک متغیر قابل آموزش، ضرب و با یک بایاس قابل آموزش، جمع می‌شود. در لایه خروجی این شبکه، برای تصمیم‌گیری از فاصله اقلیدسی تابع پایه شعاعی<sup>۳۲</sup> استفاده می‌شود و خروجی دارای کمترین فاصله از ۸۴ ویژگی لایه آخر، تصمیم شبکه را مشخص می‌کند.

نتایج به دست آمده از لینت-۱ نشان داد شبکه پیشی بزرگ‌تری برای استفاده کامل از داده ورودی موردنیاز است. به همین دلیل شبکه لینت-۴ [۱۷] با لایه‌ها و پالایه‌های بیشتر ارائه شد. این شبکه با ورودی‌های ۳۲×۳۲، دارای ۴ لایه است. لایه اول آن پیچش با اندازه ۵×۵، با ۴ نقشه ویژگی به عنوان خروجی است. در لایه دوم با ادغام دو تا دو تا از نقشه‌های لایه قبل به ۸ نقشه ویژگی دست پیدا می‌کنیم، در اینجا نیز ادغام به صورت میانگین‌گیری انجام می‌شود و تابع فعال‌سازی پس از آن اعمال می‌شود. لایه سوم پیچش با اندازه ۵×۵ به خروجی ۱۶ نقشه ویژگی، و لایه چهارم نیز ادغام به تعداد لایه قبل است، یعنی هر ماتریس لایه قبل، با میانگین‌گیری از واحدهای ۲×۲، به یک ماتریس در این لایه تبدیل می‌شود و سپس دوباره اعمال تابع فعال‌سازی انجام می‌شود. در انتهای شبکه این بار پیش از لایه تصمیم‌گیری یک لایه مخفی تمام متصل با ۱۲۰ نورون وجود دارد، و در انتها لایه تصمیم‌گیری که دارای ۱۰ نورون است، قرار دارد.

نسخه تقویت‌شده<sup>۳۱</sup> لینت-۴ از تکنیک تقویت شبکه استفاده می‌کند، که در این مورد لایه‌ها به شکل موازی آموزش داده نمی‌شوند و در هر لایه به موارد اشتباه تشخیص‌داده‌شده در لایه‌های پیشین وزن بالاتری داده می‌شود تا شبکه بتواند آن‌ها را بهتر یاد بگیرد. اولین شبکه به صورت عادی آموزش داده می‌شود، دومی ۵۰ درصد از دادگان آموزش را از میان تشخیص‌های صحیح و ۵۰ درصد را از میان تشخیص‌های غلط می‌گیرد. شبکه سوم نیز بر روی دادگانی کار می‌کند که دو شبکه قبلی در خصوص نتیجه آن توافق نداشته‌اند. علاوه بر این، داده‌هایی که در دسته‌بند ابتدایی با درصد اطمینان بالایی تشخیص داده شوند دیگر وارد مراحل بعدی نمی‌شوند. این تصمیم‌ها موجب شده با این‌که این شبکه دارای ۳ دسته‌بند است، اما میزان پردازش آن تنها ۱,۷۵ برابر لینت-۴ معمول است. به خاطر تعدد شبکه‌ها و به تبع آن، بیشتر بودن تعداد پارامترها، نیاز به دادگان ورودی بیشتری برای آموزش



شکل ۱۱: مقایسه نرخ خطای نسخه‌های مختلف معماری‌های لینت [۱۸]



شکل ۱۲: ساختار شبکه الکسنت

قبلی صورت گرفته است مانند افزایش تعداد لایه‌ها، استفاده از تابع ReLU، نرمال کردن و استفاده از ادغام بیشینه. دلیل استفاده از توابع فعال‌سازی در شبکه‌های عصبی، ایجاد قابلیت غیرخطی بودن در شبکه است که بتوان با آن تصمیمات دقیق‌تری را ایجاد کرد. در شبکه الکسنت، به دو دلیل از ReLU به‌عنوان تابع فعال‌سازی استفاده شده است و این یک تغییر اساسی است که در معماری‌های پس از الکسنت هم در بسیاری از شبکه‌ها استفاده شده است. دلیل اول اشباع‌پذیری توابع فعال‌سازی مانند سیگموئید و تانژانت هایپربولیک است و دلیل دوم ساده‌تر شدن محاسبات و افزایش سرعت آموزش شبکه است.

نرمال کردن خروجی‌های یک لایه، برای ورود به لایه بعد، به خاطر رفع مشکل اشباع توسط ReLU مورد نیاز نیست، اما استفاده از آن به تعمیم‌پذیری شبکه کمک می‌کند. این نوع از نرمال‌سازی، کار نرمال کردن را در یک نقطه، به صورت میان‌لایه‌ای، و نه در یک لایه، در همسایگی آن نقطه انجام می‌شود. هدف آن نیز تقویت مقادیر بزرگ و

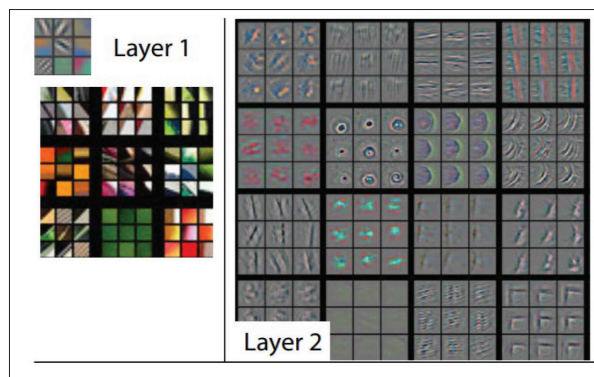
برای محاسبه فاصله، ماتریس‌های دوبعدی با اندازه  $۱۲ \times ۷$  ساخته شده است که نماینده نویسه‌های مختلف هستند (شکل ۸) و فاصله از آن‌ها محاسبه می‌شود. در آموزش این شبکه نیز برای بالا بردن قدرت تعمیم‌پذیری، تبدیل‌هایی مانند چرخش، مقیاس و جابجایی بر روی برخی ورودی‌ها اعمال می‌شود.

در مقاله [۱۸] بر روی مجموعه داده NIST [۱۹] و همچنین نسخه تغییر یافته آن<sup>۳</sup>، که با توجه به نیاز شبکه از روی داده‌های اصلی ایجاد شده‌اند، مقایسه‌ای میان انواع لینت انجام شده است که نتیجه در شکل ۹ (بر حسب مقدار خطا) قابل مشاهده است.

## ۲- الکسنت

این شبکه [۲۰] که در سال ۲۰۱۲ معرفی شده است، به دلیل کارایی بالای آن در کاربرد دسته‌بندی تصاویر، شروع توجه جدی پژوهشگران به شبکه‌های عصبی پیچشی است. در این شبکه چند تغییر نسبت به کارهای

33- Modified NIST (MNIST)



شکل ۱۳: ویژگی‌های استخراج شده توسط الکسنت در لایه اول و دوم [۲۲]

کم تاثیر کردن مقادیر کوچک است. اعمال این روش در دسته‌بندی مجموعه داده CIFAR-10، خطا را از ۱۳ درصد به ۱۱ درصد کاهش داده است.

همچنین در لایه ادغام این شبکه، برخلاف شبکه لینت از ادغام بیشینه استفاده می‌شود. ویژگی دیگری که برای کاهش بیش‌برازش در شبکه، هم‌پوشانی داشتن در بلوک‌هایی است که روی آن‌ها عمل ادغام انجام می‌شود. ساختار الکسنت در شکل ۱۰ قابل مشاهده است.

لایه اول پیش آن دارای ۹۶ پالایه  $3 \times 11 \times 11$  است. دلیل بعدی بودن پالایه‌ها، سه بعدی بودن تصویر ورودی شبکه است. عمل پیش با اندازه گام ۳ روی آن اعمال می‌شود که بر روی خروجی آن ادغام بیشینه، با اندازه ناحیه ۳ و گام ۲ اعمال شده است. سپس تابع غیرخطی و نرمال‌سازی اجرا می‌شود. لایه دوم پیش شامل ۲۵۶ پالایه  $48 \times 5 \times 5$  است که بر روی خروجی آن نیز مانند لایه قبل ادغام بیشینه با همان ویژگی‌ها اعمال شده است. دلیل عدد ۴۸ این است که ۹۶ نقشه ویژگی خروجی مرحله قبل به دو دسته ۴۸ تایی تقسیم شده است. سومین لایه پیش ۳۸۴ پالایه با اندازه  $128 \times 3 \times 3$ ، و لایه چهارم نیز ۳۸۴ پالایه با اندازه  $192 \times 3 \times 3$  را دارا است. لایه پنجم پیش، ۲۵۶ پالایه با اندازه  $192 \times 3 \times 3$  دارد که بر روی خروجی آن ادغام بیشینه با گام ۴ انجام می‌شود و خروجی آن پس از پردازش در دو لایه تمام متصل، به لایه نهایی می‌رسد. استفاده از ترکیب کردن نتیجه چندین شبکه مختلف برای تصمیم‌گیری نهایی، مانند لینت-۴ ارتقایافته، روشی

است که نتایج مثبتی را در بهبود دقت نشان داده است اما به دلیل افزایش قابل توجه میزان پردازش از آن صرف نظر شده است. ایده‌ای که الهام گرفته از ایده ترکیب شبکه‌هاست و به عنوان جایگزین مطرح شده، تکنیک حذف تصادفی<sup>۳۴</sup> [۲۱] است که شرح آن در فصل بعد آمده و در این مقاله احتمال فعال یا غیرفعال‌سازی نورون‌ها ۰/۵ در نظر گرفته شده است. با این روش، میزان پردازش به حدود دو برابر افزایش می‌یابد. این میزان پردازش بیشتر نیز صرف همگرا شدن شبکه می‌شود که بیشتر از حالت عادی به طول می‌انجامد.

شبکه الکسنت در ۲۰۱۲ با ۸ لایه در مسابقه ILSVRC و روی دادگان ImageNet [۱] با خطای حدود ۱۶ درصد و با فاصله زیاد از نفر دوم، برنده مسابقه شد و از آن به بعد توجه به استفاده از شبکه عصبی پیچشی در پردازش تصویر و توسعه این شبکه‌ها سرعت گرفته است.

### ۳- ZFnet

شبکه ZF در سال ۲۰۱۳ برنده مسابقه ILSVRC شده است که در آن پژوهشگران به تکنیکی برای بصری‌سازی ویژگی‌های استخراج شده در هر لایه از شبکه عصبی پیچشی رسیده‌اند [۲۲] و به کمک آن ویژگی‌های به دست آمده از اجرای الکسنت را به تصویر کشیده‌اند. از نتایج به دست آمده که در شکل ۱۱ قابل مشاهده است، دو مشکل در ساختار الکسنت مشاهده شده که برای حل آن‌ها تغییراتی در ساختار ایجاد شده است.

مشکل اول این‌که ویژگی‌های استخراج شده توسط الکسنت در لایه اول (شکل ۱۳: ویژگی‌های استخراج شده توسط الکسنت در لایه اول و دوم [۲۲])، بیشتر شامل تغییرات با فرکانس‌های پایین و فرکانس‌های بالا هستند و ویژگی‌هایی با فرکانس میانه کمتر در آن‌ها وجود دارد. برای رفع این مشکل اندازه پالایه لایه اول را از اندازه  $11 \times 11$  به اندازه  $7 \times 7$  تغییر داده‌اند. مشکل دوم در خصوص تداخل<sup>۳۵</sup> ویژگی‌های استخراج شده در لایه دوم است که به خاطر

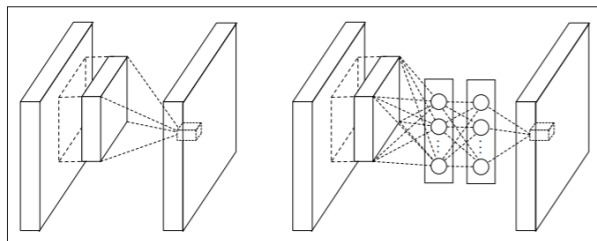
34- Drop Out  
35- Aliasing

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64	conv3-64	conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128	conv3-128	conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

شکل ۱۵: نسخه‌های مختلف معماری VGGnet [۲۶]

## ۵- VGGnet

در این معماری مربوط به سال ۲۰۱۴ [۲۶] که برنده چالش ILSVRC در همان سال شده است، طراحان تصمیم گرفته‌اند با افزایش عمق شبکه به دسته‌بندی بهتری دست پیدا کنند و در عین حال میزان پارامترها را نیز کاهش دهند. ایده این معماری شکستن پالایه‌ها به پالایه‌های کوچک‌تر است که هر دو هدف را برآورده می‌کند. در این ساختار پیچش اصلی مورد استفاده با پالایه ۳×۳ صورت می‌گیرد. در این مقاله چینش‌های متنوعی از لایه‌ها طراحی و سنجیده شده‌اند تا به یک نمونه معماری بهتر دست یابند. شش نوع معماری سنجیده شده در این مقاله در شکل ۱۳ آمده‌اند. تمامی لایه‌های پیچش ۳×۳ استفاده شده در این شبکه‌ها دارای اندازه گام یک و یک ردیف حاشیه هستند تا اندازه تصویر تا انتهای بلوک ثابت بماند. ادغام بیشینه موجود در جدول نیز با اندازه بلوک ۲×۲ با اندازه گام ۲ است، یعنی بلوک‌ها هم‌پوشانی ندارند. نکته دیگری که در این مقاله و در سنجش مدل‌های ذکر شده در بالا قابل توجه است این است که دو نوع آموزش و سه نوع سنجش برای آن در نظر گرفته شده است.



شکل ۱۴: پیچش در شبکه‌های معمولی (چپ) جایگزینی پیچش با پرسپترون چندلایه در ساختار شبکه در شبکه (راست) [۲۳]

اندازه گام بزرگ پیچش رخ داده، به همین دلیل اندازه گام از ۴ به ۲ تغییر یافته است. در نتیجه این تغییرات ویژگی‌های متنوع‌تر و متمایزتری با روش ارائه شده به دست آمده‌اند و کارایی شبکه را در مقایسه با الکسنت بهبود داده‌اند.

## ۴- شبکه در شبکه

ایده این پژوهش [۲۳]، مبنای شکل‌گیری ایده استفاده از پیچش ۱×۱ است، که در شبکه‌هایی مانند VGG، گوگل‌نت [۲۴] و رزنت ۳۷ [۲۵]، برای تغییر تعداد لایه‌ها و عموماً کاهش میزان پردازش مورد استفاده قرار گرفته است. این پژوهش استفاده از ساختار پرسپترون چندلایه را به جای استفاده از لایه پیچش پیشنهاد می‌دهد و دلیل آن را این اعلام می‌کند که شبکه عصبی پیچشی یک عملگر خطی است که ما پس از انجام آن یک عملیات جهت غیرخطی‌سازی انجام می‌دهیم، در صورتی که استفاده از یک تخمین‌گر جهانی توابع ۳۸ مانند پرسپترون چندلایه که توانایی بالاتری برای ساخت توابع غیرخطی دارد، کارایی شبکه را افزایش می‌دهد. نحوه عملکرد نیز بدین صورت است که پرسپترون چندلایه، به عنوان یک میکروشبکه، جایگزین عملگر پیچش می‌شود و ناحیه ورودی پیچش، ورودی آن قرار می‌گیرد. مثال ارائه شده در مقاله، میکروشبکه‌ای سه لایه است که در لایه‌های آن از ReLU برای غیرخطی‌سازی استفاده می‌شود (شکل ۱۴).

همچنین ایده دیگری که در این مقاله مطرح می‌شود، استفاده از ادغام کلی میانگین [۳۶] است، که برای جایگزینی لایه‌های تمام متصل در انتهای شبکه‌ها استفاده شده است.

- 36- GoogleNet
- 37- ResNet
- 38- Universal Function Approximator
- 39- Global Average Pooling

انواع آموزش اعمال شده، تک مقیاس<sup>۴۰</sup> و چند مقیاس<sup>۴۱</sup> هستند، آموزش تک مقیاس مشابه حالت عادی آموزش است و آموزش چند مقیاسی، با این هدف انجام می‌شود که یک جسم خاص در تصاویر مختلف، برحسب فاصله می‌تواند در اندازه‌های متفاوتی ظاهر شود، به همین دلیل نیاز است شبکه بتواند در اندازه‌های مختلف آن را تشخیص دهد. بدین خاطر اندازه تصاویر ورودی را به یک مقدار در بازه [کمینه: بیشینه] تغییر داده و سپس بر روی قطعه‌ای از آن به اندازه ورودی شبکه، مرحله آموزش را انجام می‌دهند.

در خصوص ارزیابی و استفاده، سه روش مطرح شده تک‌برش<sup>۴۲</sup>، چندبرش<sup>۴۳</sup> و ارزیابی چگال<sup>۴۴</sup> است. استفاده از روش چندبرش برای ارزیابی، بدین صورت انجام می‌شود که بر روی برش‌های متفاوتی از ورودی و نسخه‌های دوران‌یافته آن‌ها ارزیابی صورت می‌پذیرد و تصمیم نهایی از میانگین نتایج قطعات به دست می‌آید. روش چگال بدین صورت است که شبکه آموزش‌دیده را نه بر روی برش‌هایی از تصویر که بر روی کل تصویر اعمال می‌کنند، و لایه‌های تمام متصل انتهای شبکه را نیز به لایه‌های پیچش تبدیل می‌نمایند. در نهایت برای هر تصویر به جای یک خروجی و یک تصمیم، چند تصمیم خواهیم داشت که با میانگین و بیشینه گرفتن بر روی نتایج، نتیجه نهایی به دست می‌آید. یک مزیت این روش نسبت به برش دادن تصویر ورودی جلوگیری از محاسبه‌ی چندباره‌ی شبکه بر روی برخی پیکسل‌ها است، ضمن این‌که در واقع تمام حالات برش‌های مختلف را نیز پوشش می‌دهیم، نه یک تعداد معین از آن‌ها را. مزیت دیگر آن این است که در حالت برش خورده، لایه پیچش در اطراف تصویر حاشیه صفری که ما داده‌ایم را برای محاسبات استفاده می‌کند اما در اینجا مقادیر واقعی در آن دخیل هستند و در واقع اصل داده بیشتر در تصمیم‌گیری موثر است.

- 40- Single Scale  
41- Multi Scale  
42- Single Crop  
43- Multi Crop  
44- Dense Evaluation

## ۶- گوگل‌نت

ساده‌ترین راه برای بهبود دقت شبکه‌های عصبی پیچشی، افزایش عمق و عرض آن‌هاست که افزایش عمق، به معنای افزایش تعداد لایه‌ها و افزایش عرض، به معنای بالا بردن تعداد ویژگی‌ها در هر لایه است. که این می‌تواند یک روش سریع -از منظر طراحی- برای بهبود دقت باشد. اما این کار دو مشکل اصلی دارد. یکی افزایش تعداد پارامترهاست که می‌تواند به بیش‌برازش در شبکه بیانجامد. خصوصاً در صورتی که تعداد نمونه‌های برچسب‌دار در مجموعه دادگان محدود باشد، این موضوع تبدیل به مشکلی جدی در یادگیری شبکه می‌شود. مشکل دیگری که افزایش مداوم عمق شبکه‌ها به دنبال دارد، افزایش نیاز به منابع پردازشی است. مثلاً در یک شبکه عصبی پیچشی، هر تعداد افزایش در تعداد فیلترهای دو لایه پیچش پشت سر هم، به صورت درجه‌دو منجر به افزایش میزان پردازش می‌شود. که اگر ساختار اضافه شده کارایی بالایی نداشته باشد و منجر به ایجاد ماتریس‌های تنک<sup>۴۵</sup> با درایه‌های متعدد نزدیک به صفر شود، هدر رفتن منابع پردازشی را به همراه خواهد داشت.

راه اصلی برای رفع هر دو مشکل می‌تواند استفاده بیشتر از ساختارهای تنک، به جای ساختارهای تمام متصل باشد، در این راستا لایه‌های تمام متصل انتهای شبکه‌های عصبی پیچشی رایج حذف شده و تنها یک لایه باقی می‌ماند. با توجه به این‌که پردازنده‌ها ساختار مناسب برای استفاده از این ساختارهای تنک را ندارند، با تبدیل آن‌ها به عملیات چگال انجام می‌شود. برای همین پژوهشگران این مقاله در سال ۲۰۱۵ [۲۴]، به دنبال ساختارهایی هستند که عملیات تنک را به کمک ساختارهای چگال پیاده‌سازی کنند. در واقع به دنبال بلوکی هستند با ویژگی مذکور طراحی شود و قابلیت تکرار شدن پشت‌سرهم برای ساختن یک شبکه عمیق را در کنار استفاده بهینه از منابع پردازشی داشته باشد. حالت پایه بلوک ارائه‌شده در شکل ۱۴ قابل مشاهده است.

در یک شبکه عصبی پیچشی معمولی، پالایه‌ها ثابت و

45- Sparse

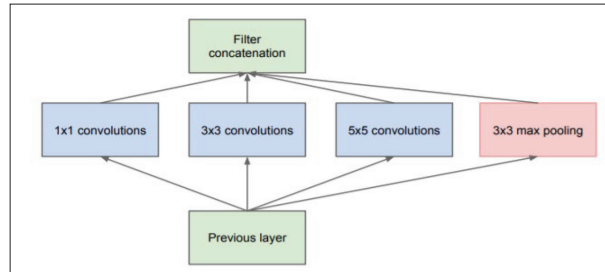
که طبق شکل ۱۶ پس از بلوک‌های دوم، هفتم و نهم لایه ادغام اجباری به صورت مجزا از لایه ادغام موازی که درون ماژول‌ها موجود است، قرار گرفته است. مورد دیگر استفاده از ادغام کلی میانگین به جای لایه‌های متصل میانی است که افزایش دقت را به دنبال دارد.

با توجه به عمیق بودن این ساختار، انتشار مؤثر گرادینت‌ها از لایه آخر به اولین لایه یک مسئله به حساب می‌آید. طبق پژوهش [۲۴] از دسته‌بندی‌های میانی برای بهبود آموزش شبکه استفاده شده است.

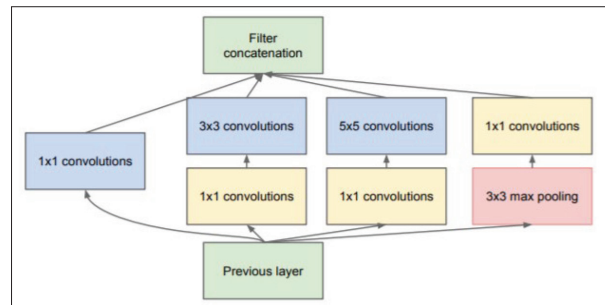
نسخه بعدی معماری گوگل با نام اینسپشن-۲ در مقاله [۲۷] ارائه شده است. در این مقاله ابتدا اصولی برای طراحی شبکه‌ها ذکر شده و سپس معماری مورد نظر ارائه شده است، این اصول عبارتند از: ۱. شبکه نباید دچار کمبود اطلاعات ورودی شود، بدین معنا که کاهش اندازه ورودی، خصوصاً در مراحل ابتدایی نباید به یکباره انجام شود، بلکه این کاهش باید در طول شبکه و به صورت تدریجی صورت پذیرد. ۲. افزایش تعداد نقشه‌های ویژگی در یک لایه کمک می‌کند ویژگی‌های به دست آمده وابستگی کمتری به یکدیگر داشته باشند که این موضوع، سرعت آموزش ویژگی‌ها را افزایش می‌دهد. ۳. کاهش عمق ورودی پیش از اعمال لایه‌ها، مشکل جدی در بازنمایی دادگان ایجاد نمی‌کند و موجب کاهش دقت زیادی نمی‌شود. ۴. برای دستیابی به شبکه‌ای بهینه، نیاز است تناسبی از عمق و عرض شبکه - به معنای تعداد پالایه‌ها در هر لایه - را در نظر بگیریم. بر اساس یک قاعده کلی انتظار می‌رود افزایش عمق و عرض شبکه ما را به شبکه‌ای با کیفیت‌تر برساند، اما با توجه به محدود بودن توان پردازشی، باید به حالت بهینه‌ای از عرض و عمق متناسب با آن دست یافت.

در معماری ارائه شده چند ایده برای بهبود ساختار به کار گرفته شده است. ایده اول، همان شکستن پیچش‌ها به پیچش‌های کوچکتر است که هم به شکل متقارن و هم نامتقارن مورد استفاده قرار گرفته است.

ایده دیگری که به کار گرفته شده است، کاهش اندازه



شکل ۱۶: بلوک پایه معماری گوگل نت [۲۴]

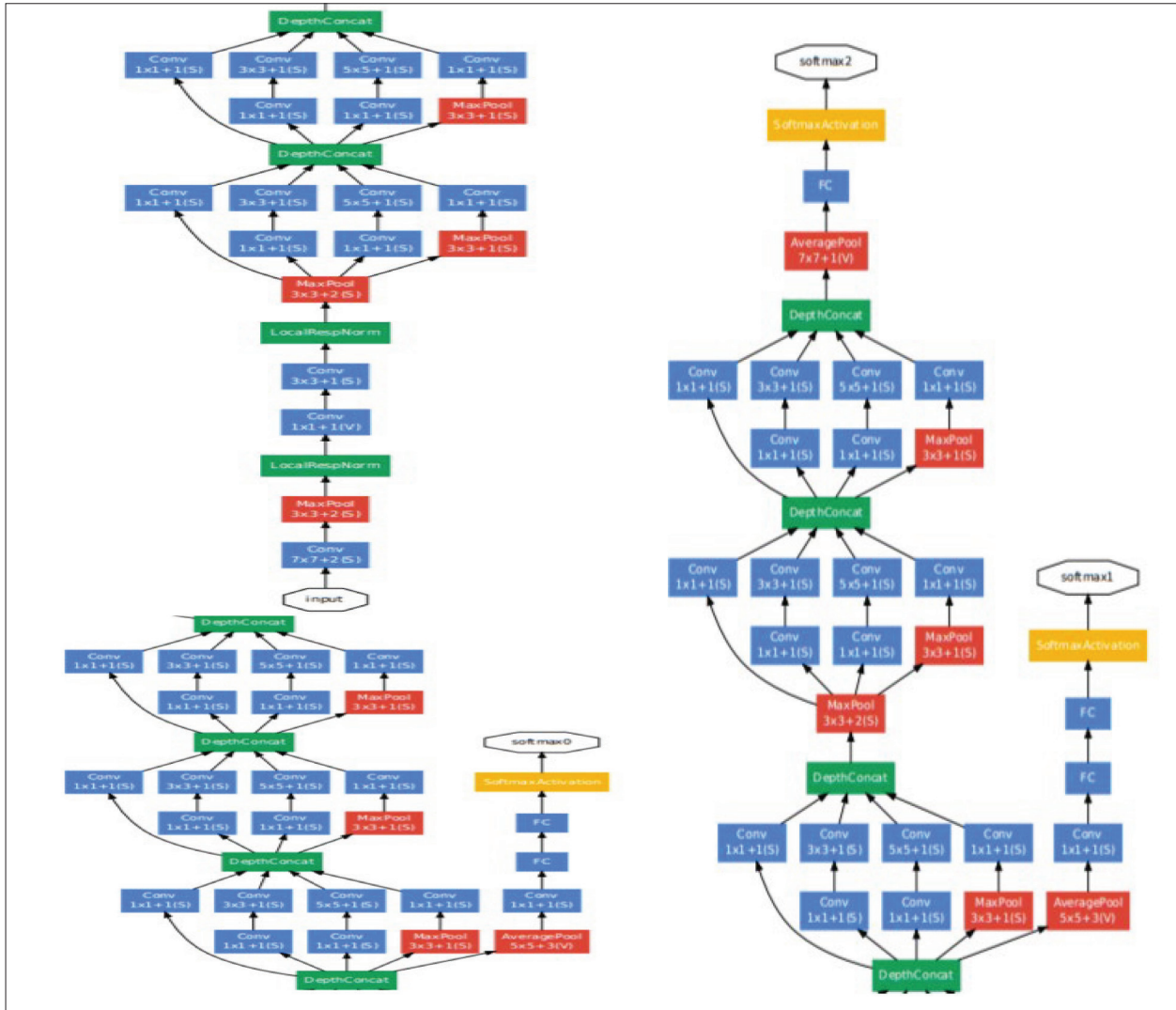


شکل ۱۷: نسخه دوم از ماژول پایه معماری گوگل نت، که به نام طراحی گلوگاهی شناخته می‌شود [۲۴]

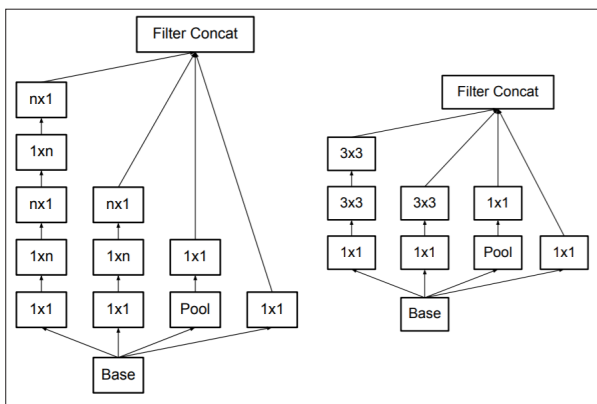
از پیش انتخاب شده هستند، اما در شبکه گوگل، در داخل هر بلاک که ماژول اینسپشن<sup>۶۱</sup> نام دارد سه عمل پیچش با اندازه پالایه‌های  $1 \times 1$ ،  $3 \times 3$  و  $5 \times 5$  را در کنار هم داریم، و در کنار آن‌ها لایه ادغام نیز وجود دارد. لایه ادغام موجود مانند انجام عمل ادغام بر روی خروجی لایه قبل است. خروجی تمامی این عملیات به صورت چگال در کنار هم قرار گرفته خروجی بلوک را می‌سازند. شبکه لایه بعد می‌تواند از هر کدام از این بخش‌ها استفاده نماید و این بدان معناست که شبکه آموزش داده شده لایه‌های فعال ثابت ندارد و بر اساس نیاز دادگان و دسته‌بندی‌های آن‌ها، از خروجی انواع پیچش استفاده می‌شود.

در نسخه بعدی این شبکه از ماژول اینسپشن که در شکل ۱۵ تصویر شده است، برای کاهش تعداد پارامترها از پیچش  $1 \times 1$  استفاده شده است. این ماژول با حفظ ابعاد طول و عرض، تعداد لایه‌های به دست آمده از پیچش را، کاهش می‌دهد.

ابتدای ساختار شبکه گوگل، مانند ساختارهای معمول شبکه عصبی پیچشی دارای دو لایه پیچش و دو لایه ادغام در میان آن‌هاست و پس از آن ۹ ماژول اینسپشن قرار دارد



شکل ۱۸: ساختار کلی گوگل نت - نسخه یک [۲۴]

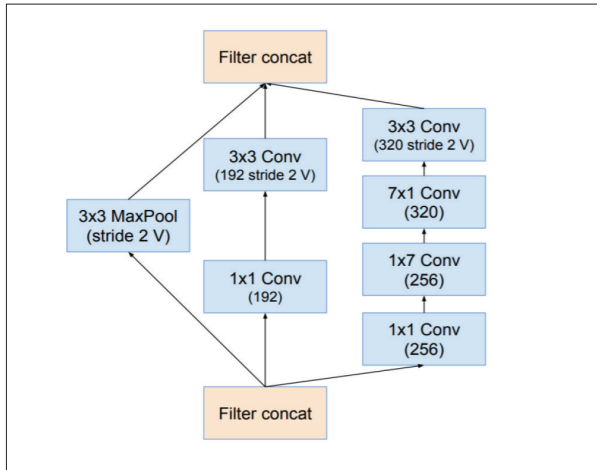


شکل ۱۹: ماژول پایه با پیچش های  $3 \times 3$  (سمت راست) استفاده از پیچش های  $1 \times n$  و  $n \times 1$  به جای  $n \times n$  پایه (سمت چپ) [۲۷] آن‌ها در استفاده از برخی از ایده‌های بالا در ماژول‌های سازنده است. نسخه معماری که تمامی ایده‌های مطرح شده

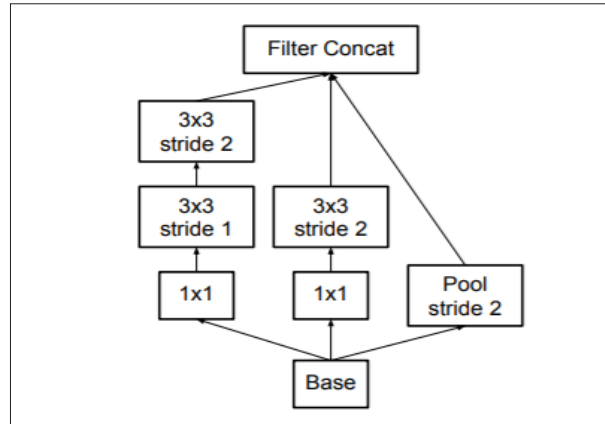
نقشه‌های ویژگی به‌گونه‌ای است که میزان پردازش مورد نیاز در عمل پیچش را کاهش دهد. استفاده از پیچش و سپس ادغام، یک چپش رایج در شبکه‌های عصبی پیچشی است، با فرض این‌که ادغام انجام شده  $2 \times 2$  باشد عملاً ۷۵٪ پردازش انجام شده بی‌استفاده باقی می‌ماند و توان پردازشی هدر می‌رود. همچنین قرار دادن ادغام پیش از پیچش، با اصل یک تناقض دارد چون حجم قابل توجهی از دادگان را حذف می‌کند و یک گلوگاه اطلاعاتی حساب می‌شود. راه‌حل ارائه‌شده استفاده از پیچش با اندازه گام ۲ به جای این دو لایه پشت سر هم است (شکل ۱۸).

معماری اینسپشن-۲ دارای سه نسخه است که تفاوت





شکل ۲۲: بلوک کاهش از اندازه  $17 \times 17$  به  $8 \times 8$  [۲۸]



شکل ۲۰: استفاده از پیچش با گام ۲ به جای ادغام [۲۷]

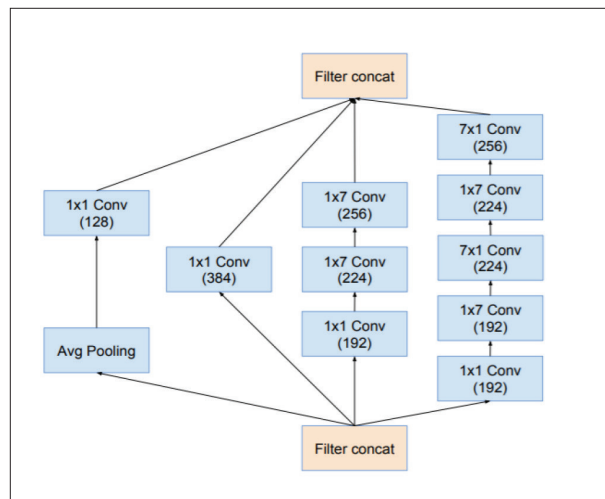
ورودی دستخوش تغییر نمی‌شود، برای کاهش اندازه ورودی بلوک‌های دیگری تحت عنوان بلوک‌های کاهش<sup>۴۸</sup> طراحی شده‌اند که ایده مطرح شده در مقاله [۲۷] در خصوص افزایش اندازه گام پیچش در آن به کار رفته است. در شکل ۱۹ و ۲۰ دو نمونه از ماژول‌های گفته شده قابل مشاهده هستند.

در معماری اینسپیشن-رزنت [۲۸] نیز بلوک‌های اینسپیشن، با میان‌برهای باقیمانده ترکیب شده و بلوک‌های پایه را مانند شکل ۲۱ در آورده‌اند. همان‌طور که قابل مشاهده است، پیش از انجام عمل جمع ورودی، به کمک یک پیچش  $1 \times 1$  تعداد لایه‌های منتج از سایر عملیات انجام شده در بلوک را با ابعاد ورودی یکسان‌سازی می‌کنیم تا عمل جمع به صورت نظیر به نظیر قابل انجام باشد.

### ۷- رزنت

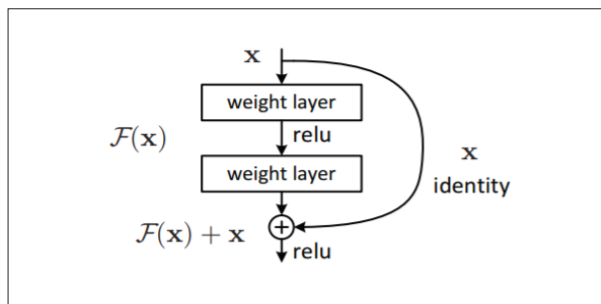
در سال‌های اخیر، در حوزه شبکه‌های عصبی پیچشی، افزایش عمق به عنوان راهکاری برای افزایش دقت این شبکه‌ها در نظر گرفته می‌شود. اما مسئله‌ای که در خصوص افزایش عمق ما را دچار مشکل می‌سازد، کوچک شدن گرادینت‌ها و به روز نشدن آن‌ها طی مراحل آموزش است، که به آن محو گرادینت<sup>۴۹</sup> می‌گویند.

در مقاله [۲۵] مربوط به سال ۲۰۱۶ که ارائه‌دهنده



شکل ۲۱: بلوک اینسپیشن B مورد استفاده در ساختار اینسپیشن-۴ [۲۸]

در این مقاله را در خود دارد، اینسپیشن-۳ [۲۷] است. مقاله دیگری که در آن نسخه‌های بعدی این معماری ارائه می‌شود، [۲۸] است که نسخه چهارم اینسپیشن و نسخه‌های بعدی که در ترکیب با معماری رزنت ارائه شده‌اند را معرفی می‌کند. در نسخه چهارم، تلاش شده است شبکه یکدست‌تری ساخته شود. در نتیجه پایه شبکه ۷×۷، یعنی بخش ابتدایی آن، دچار تغییر شده است و علاوه بر آن، سه نوع بلوک اینسپیشن A، B، و C برای ورودی با اندازه‌های  $35 \times 35$ ،  $17 \times 17$  و  $8 \times 8$  طراحی و در ساختار مورد استفاده قرار گرفته است. در هر بلوک پیچش‌ها با اندازه‌های متفاوتی قرار دارند که تعداد پالایه‌های هر کدام به صورتی انتخاب شده است که توان پردازشی به صورت متعادل در شاخه‌ها مورد استفاده قرار گیرد. در این بلوک‌های اینسپیشن، اندازه



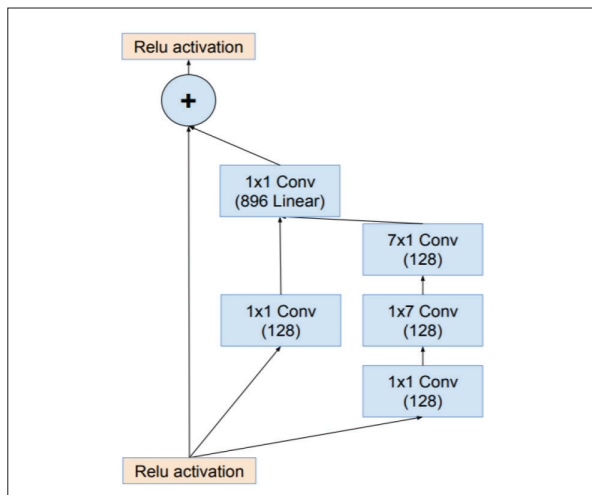
شکل ۲۴: ساختار پایه بلوک رزنت [۲۵]

موارد، تعداد لایه‌های نقشه‌ی ویژگی‌ها را افزایش می‌دهند. در انتهای این شبکه نیز یک لایه ادغام کلی میانگین و یک لایه تمام متصل برای تصمیم‌گیری، به همراه softmax وجود دارد.

در این مقاله [۲۸] نوع دیگری از ساختار رزنت و سه نمونه از آن نیز ارائه می‌شود. این نوع از ساختار را، ساختار گلوگاهی<sup>۱۱</sup> نامیده‌اند. ایده آن این است که در پیچش‌های  $3 \times 3$  که پرهزینه‌ترین فعالیت ساختار رزنت به حساب می‌آید، تعداد نقشه‌های ویژگی را به کمک پیچش  $1 \times 1$  کاهش داده و پس از انجام پیچش  $3 \times 3$ ، آن را به تعداد ابعاد موردنیاز بازگردانند. شکل ۲۴ نشان‌دهنده بلوک‌های ساختار گلوگاهی شامل پیچش  $1 \times 1$  برای کاهش ابعاد، پیچش  $3 \times 3$  برای استخراج ویژگی و یک پیچش  $1 \times 1$  دیگر برای بازگرداندن ابعاد به تعداد موردنظر است. اگر بخواهیم همان ساختار ۱۶ بلوکی (۳۴ لایه) را با بلوک‌های جدید بسازیم، دارای ۵۰ لایه خواهد بود.

در این پژوهش، برای سنجش نتایج، یک ساختار ۳۴ لایه بدون اتصال‌های میان‌بر، ساختار ۳۴ لایه اولیه رزنت و ساختار ۵۰ لایه گلوگاهی را بر روی مجموعه دادگان CIFAR-10 اعمال کرده‌اند که بهبود دقت را به میزان ۰,۵۴ درصد نشان داده است.

همچنین دو نمونه دیگر از ساختار گلوگاهی با تعداد لایه‌های ۱۰۱ و ۱۵۲، برای دسته‌بندی مجموعه دادگان ImageNet که حجم بالاتری دارد ساخته شده است. تا اینجا با افزایش بلوک‌ها شاهد افزایش دقت بوده‌ایم. اما نمونه‌ای از ساختار گلوگاهی با ۱۲۰۲ لایه ساخته شده است که

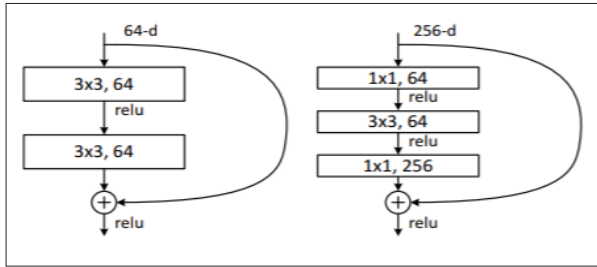


شکل ۲۳: بلوک اینسپشن B ترکیب‌شده با اتصال میان‌بر [۲۸]

معماری رزنت است، بیان شده است که برای رفع این مشکل، می‌توان از ساختارهای مبتنی بر باقیمانده<sup>۱۲</sup> استفاده کرد، و این فرضیه بر مبنای بهبودهایی است که در استفاده از روش‌های مبتنی بر باقیمانده در مسائل دیگر و در روش‌های مبتنی بر ایجاد میان‌بر در شبکه‌ها دیده شده است. راهکار ارائه شده برای رفع این مشکل در این پژوهش، استفاده از اتصال میان‌بر و ساختن شبکه‌ای مبتنی بر بازنمایی باقیمانده است.

ساختار اولیه این شبکه شامل ۳۴ لایه است و ساختار آن با الهام از VGGnet، به جز در لایه اول، متشکل از پیچش‌هایی با اندازه  $3 \times 3$  است. این شبکه شامل ۱۶ بلوک است. بلوک اولیه سازنده شبکه، شامل یک جفت پیچش  $3 \times 3$ ، با یک لایه ReLU - برای غیرخطی کردن - در میان آن‌ها و اتصال میان‌بر است که پس از جمع شدن مقدار ورودی با خروجی پیچش‌ها، بار دیگر ReLU بر آن اعمال می‌شود (شکل ۲۲). نکته‌ای که در خصوص این ساختار وجود دارد این است که اتصال‌های میان‌بر، پارامتر جدیدی به شبکه اضافه نمی‌کنند، که این موضوع، امکان مقایسه با شبکه‌های ساده با عمق مشابه را امکان‌پذیر می‌سازد.

ساختار کلی شبکه رزنت در شکل ۲۳ نشان داده شده است. همان‌طور که در این شکل دیده می‌شود، برخی از اتصال‌های میان‌بر به صورت خط چین ترسیم شده‌اند که این



شکل ۲۶: ماژول اولیه رزنت (چپ) و ماژول کلوگاهی رزنت (راست) [۲۵]

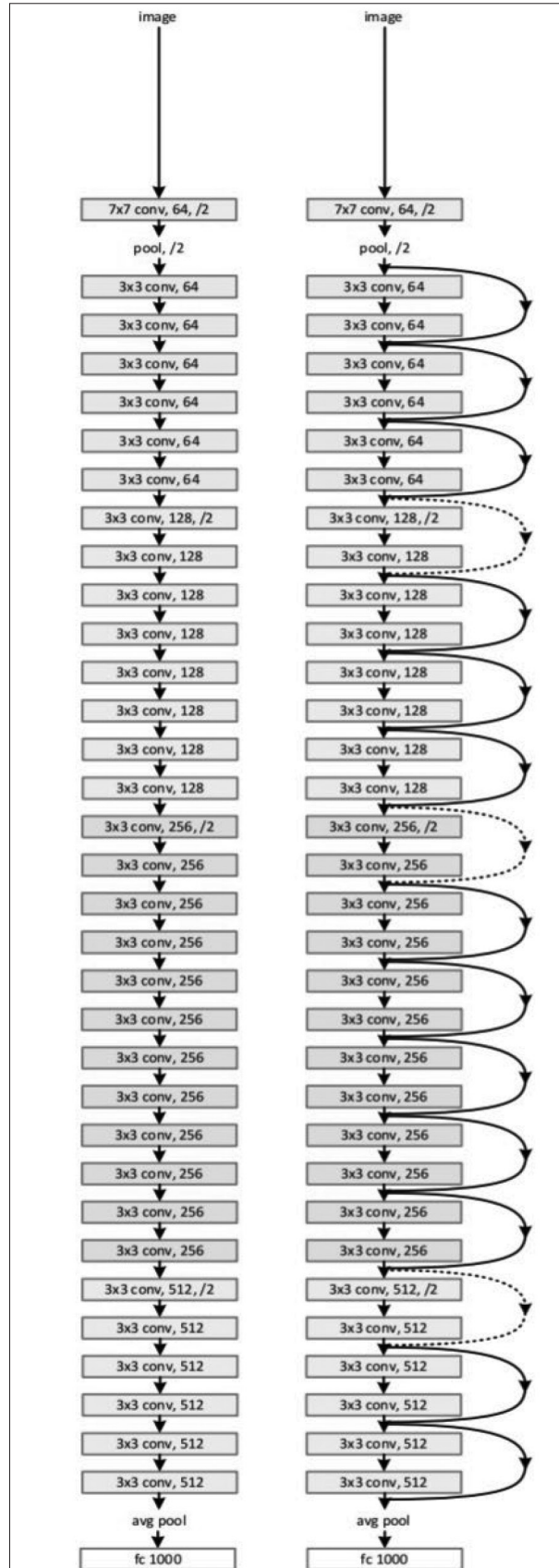
میزان دقت آن از ساختار ۱۵۲ لایه کمتر بوده است، و حدس پژوهشگران مقاله بر این است که این افزایش بالای تعداد پارامترهای شبکه، موجب بیش‌برازش شده است.

نکته‌ای که در خصوص استفاده از میان‌برها در شبکه وجود دارد این است که خروجی هر بلوک باید ابعاد مساوی با ورودی آن داشته باشد تا بتوان جمع نظیر به نظیر را در آن انجام داد. در صورتی که تعداد لایه‌های ورودی و خروجی برابر نباشند، با ضرب یک ماتریس تصویر خطی<sup>۲</sup> از ماتریس ورودی را، با ابعاد موردنیاز می‌سازیم.

#### ۷-۱- ایده‌های بهبود ساختار رزنت

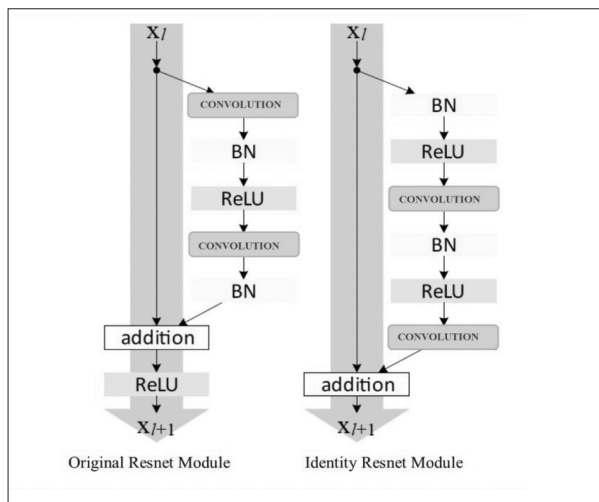
در سال‌های اخیر اکثر پژوهش‌ها به بهبود ساختارهای مطرح‌شده مانند رزنت و گوگل‌نت می‌پردازند و کمتر مواردی یافت می‌شود که ساختار را از پایه به نوع دیگری می‌سازند. برخی از این موارد در زیر ادامه بررسی می‌شوند.

در ساختار رزنت، خروجی هر بلوک، ترکیب مشخصی از ورودی و نتیجه عملیات درون بلوک است. ایده مطرح‌شده در مقاله [۲۹] که هایوینت<sup>۳</sup> نام‌گذاری شده است، پارامتر دیگری را به هر بلوک می‌افزاید و آن ضربی است که مشخص می‌کند چه نسبتی از خروجی متعلق به ورودی و چه نسبتی متعلق به سایر عملیات درون بلوک باشد. این ایده که الهام گرفته از معماری شبکه عصبی حافظه کوتاه‌مدت ماندگار<sup>۴</sup> است، در واقع یک سازوکار درجه‌ای<sup>۵</sup> پیش از خروجی بلوک قرار می‌دهد (فرمول ۳).



شکل ۲۵: ساختار یک شبکه معمولی ۳۴ لایه (چپ) در مقایسه با رزنت ۳۴ لایه (راست) [۲۵]

52- Linear Projection  
53- HighwayNet  
54- Long Short-Term Memory (LSTM)  
55- Gating Mechanism



شکل ۲۷: ساختار ماژول رزنت تغییر یافته در مقاله [۳۰]

جدول ۱: دقت‌های به دست آمده از سنجش انواع مختلف پیش درون بلوک رزنت [۳۱]

block type	depth	# params	time,s	CIFAR-10
$B(1, 3, 1)$	40	1.4M	85.8	6.06
$B(3, 1)$	40	1.2M	67.5	5.78
$B(1, 3)$	40	1.3M	72.2	6.42
$B(3, 1, 1)$	40	1.3M	82.2	5.86
$B(3, 3)$	28	1.5M	67.5	5.73
$B(3, 1, 3)$	22	1.1M	59.9	5.78

شبکه‌هایی با عمق ۱۶، ۲۲، ۲۸ و ۴۰ آزمایش شده است. در شبکه‌هایی با عمق ۱۶ و ۲۲ و ۲۸ با افزایش ضریب دقت کاهش یافته است، اما در شبکه با عمق ۴۰، با افزایش عرض با ضریب ۸، دقتی کمتر از شبکه ۲۲ لایه با ضریب ۸ کسب شده است. در مقایسه شبکه‌های این آزمایش با شبکه‌های رزنت عادی، می‌توان اشاره داشت به شبکه ۴۰ لایه با عرض ۴ برابر، که توانسته است به دقتی بالاتر از رزنت ۱۰۰۱ لایه به دست آورد در حالی که میزان پارامترهای آن نیز از شبکه هزار لایه کمتر بوده است. این تغییر میزان پارامترهای شبکه را با ضریب مشخص شده افزایش می‌دهد اما زمان آموزش و بهینه شدن پارامترها نسبت به شبکه‌هایی با عمق بیشتر، به دلیل امکان پردازش موازی بسیار کمتر است.

مقاله بعدی [۳۲]، در تلاش برای آموزش شبکه‌های باقیمانده بسیار عمیق، از ایده حذف تصادفی برای بلوک‌های رزنت استفاده کرده است که در هر دور از

$$y_i = H_i(X) * T_i(X) + x_i * (1 - T_i(X)) \quad (۳)$$

این سازوکار برای بهینه‌سازی شبکه‌های عمیق ارائه شده است، بدین صورت که سیگنال را به صورت تطبیق پذیر برای عمق مربوطه در شبکه عبور می‌دهد، که این امر به جریان بهتر گرادیان در مرحله پس‌انتشار کمک می‌کند.

در مقاله دیگری [۳۰] انواع تغییر یافته‌ای از بلوک رزنت مورد بررسی قرار گرفته و نهایتاً تغییر در چینش تبدیل‌ها ارائه شده است که به جریان بهتر گرادیان‌ها در آن کمک می‌کند. این چینش جدید امکان آموزش شبکه‌های عمیق‌تر از ۱۰۰۰ لایه را پدید می‌آورد. به بلوک ساخته شده در این مقاله رزنت با پیش‌فعال‌سازی<sup>۵۶</sup> گفته می‌شود، که در آن عملیات نرمال‌سازی دسته‌ای و ReLU پیش از عمل پیش انجام می‌شوند، این عملیات نرمال‌سازی تأثیری مانند منظم‌سازی<sup>۵۷</sup> بر روی داده‌ها دارد. همچنین در این مقاله سعی شده است کمترین میزان تغییر بر روی خروجی بلوک صورت پذیرد، به همین دلیل عملیات ReLU پس از جمع شدن دو مقدار خروجی از بلوک حذف شده و در شاخه‌ی مربوط به عمل پیش قرار گرفته است (شکل ۲۵). مقاله دیگری با عنوان شبکه‌های باقیمانده عریض<sup>۵۸</sup> [۳۱]، تلاش کرده است با تغییر متغیرهای معمول در بلوک‌های رزنت مانند عمق، عرض و نوع پیش‌مورد استفاده تأثیر این عوامل را در بهبود کارایی بررسی کند. موارد زیر برخی از دستاوردهای این مقاله هستند:

a. از نظر عمق بلوک، یا همان تعداد پیش‌های درون بلوک، دو پیش بهترین عملکرد را داشته و با افزایش آن، تعداد اتصال‌های میان‌بر در شبکه کاهش یافته و به دلیل مشکل شدن بهینه‌سازی، دقت شبکه کاهش می‌یابد.

b. از نظر نوع پیش‌های داخل بلوک، حالات مختلفی در نظر گرفته شده است مانند که حالت ساده یعنی (۳×۳) که بلوک شامل دو پیش با اندازه ۳ در ۳ است به بالاترین دقت رسیده است (جدول ۱).

c. افزایش عرض بلوک، با ضرایب ۱ الی ۱۲ بر روی

56- Pre-Activation  
57- Regularization  
58- Wide Residual Networks

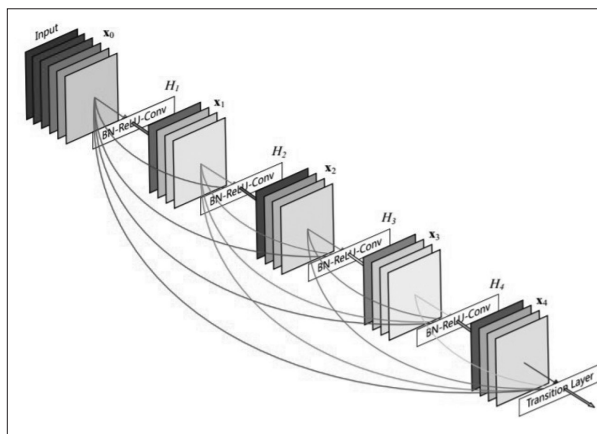
می‌شود. این موضوع نویسندگان این مقاله را به این نتیجه رسانده است که عمق خیلی زیاد ایده‌آل نیست و احتمال افزونگی در شبکه‌های بسیار عمیق بالاست. به همین دلیل سعی شده است به کمک باز استفاده از ویژگی‌ها، شبکه‌ای با دقت بالاتر ساخته شود.

بلوک چگال ۶۰ که در این مقاله طراحی می‌شود، دارای چند مرحله عملیات نرمال‌سازی - غیرخطی سازی و پیچش است که در هر مرحله تمامی نقشه‌های ویژگی مراحل پیشین به‌عنوان ورودی در نظر گرفته می‌شوند، و در واقع دانش شبکه به‌صورت تجمعی در اختیار آن مرحله قرار می‌گیرد، بدین‌صورت شبکه در هر لایه می‌تواند ویژگی‌هایی از طیف ساده تا پیچیده را برحسب نیاز حدس بزند. یکی از پارامترهای مهم در این شبکه نرخ رشد است که در واقع همان تعداد نقشه‌های ویژگی در هر مرحله است که به موارد پیشین خود افزوده می‌شود. نرخ رشد عموماً عدد بزرگی نیست و این نکته‌ای است که موجب کاهش میزان پارامترهای این معماری نسبت به سایر معماری‌ها می‌شود. اعداد پیشنهاد شده برای نرخ رشد اعدادی مابین ۱۲ و ۴۰ هستند که در مقایسه با تعداد نقشه‌های ویژگی در هر لایه از سایر شبکه‌ها عدد کوچکی است. در میان بلوک‌های چگال این شبکه لایه‌های گذار قرار دارند که شامل پیچش  $1 \times 1$  و لایه ادغام با اندازه  $2 \times 2$  هستند.

#### ۹- اسکوییزنت<sup>۶۱</sup>

در نوعی دیگر از مقالات این حوزه، به جای تلاش برای دستیابی به دقت بالاتر، تلاش می‌شود میزان پارامترها و حجم مدل کاهش یابد [۳۵] it is typically possible to identify multiple DNN architectures that achieve that accuracy level. With equivalent accuracy, smaller DNN architectures offer at least three advantages: (۱) برخی مزایا که برای ساخت مدل‌های کوچک ذکر شده است عبارتند از: سرعت بالاتر و سربار کمتر در آموزش مدل روی کارسازهای موازی، امکان انتقال سریع‌تر به روزرسانی‌ها به کامپیوترهای کاربران و امکان

60- Dense Block  
61- SqueezeNet



شکل ۲۸: ساختار درون یک بلوک چگال با نرخ رشد ۴ [۳۴]

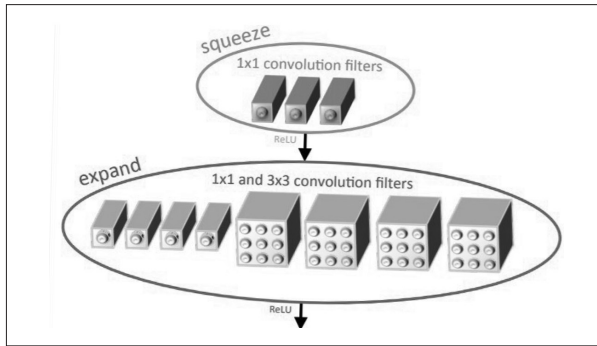
آموزش، تعدادی از بلوک‌های رزنت را به‌صورت تصادفی غیرفعال می‌کند و ورودی بلوک از میان‌بر به‌صورت بدون تغییر عبور می‌کند. البته این صرفاً مربوط به زمان آموزش است، و در زمان سنجش تمام شبکه مورد استفاده قرار می‌گیرد.

پژوهش دیگری با ارائه ساختاری به نام ResNext [۳۳]، تلاش کرده است با افزایش عرض شبکه بدون تغییر میزان پارامترها به دقت بالاتری دست یابد. در اینجا متغیر دیگری با عنوان قدرتمندی ۵۹ مطرح می‌شود، که هر بلوک از رزنت را تبدیل به بلوکی با تعدادی مسیر موازی می‌کند که هر یک از این مسیرها به میزان عدد قدرتمندی دارای نقشه ویژگی هستند. این مقاله با اعمال تابعی برای غیرخطی سازی در هرکدام از این مسیرها، سعی در افزایش کارایی آن‌ها دارد. این موضوع موجب افزایش دقت شبکه ارائه شده می‌گردد.

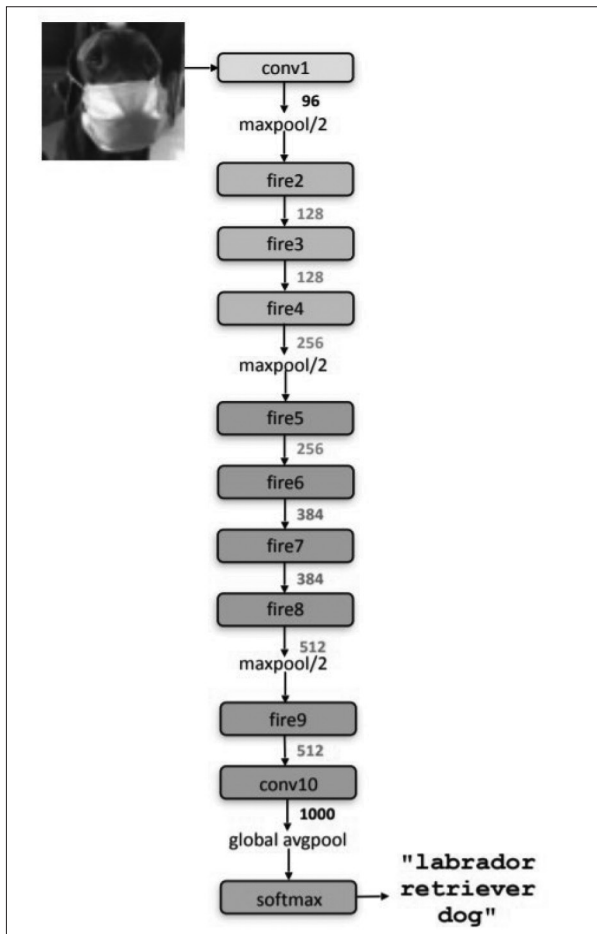
#### ۸- دِنسِنِت

این شبکه [۳۴]، با ایده‌ای شبیه به رزنت اما متفاوت از آن ساختاری را طراحی کرده است که ویژگی بارز آن باز استفاده از ویژگی‌های ساخته شده در لایه‌های پیشین است. ایده این مقاله برگرفته از مشاهداتی است که در مقالات پس از معماری رزنت در تلاش برای بهینه‌سازی شبکه‌های عمیق‌تر از ۱۰۰ یا هزار لایه دیده شده است. مقالاتی مانند [۲۹] و [۳۲]، که تغییری که در شبکه - برای آموزش آن - ایجاد می‌کنند موجب گذر و نادیده گرفته شدن برخی لایه‌ها

59- Cardinality



شکل ۲۹: ساختار ماژول آتش، ریزساختار سازنده اسکویزنت [۳۵]



شکل ۳۰: ساختار اسکویزنت [۳۵]

دقت مشابه الکسنت دست پیدا می‌کنند. حجم مدل ساخته شده حدود  $4/8$  مگابایت است که با این وجود تلاش شده است با استفاده از روش‌های فشرده‌سازی به حجم پایین‌تری دست یابند. حجم کوچکترین مدل به دست آمده در این پژوهش  $0/47$  مگابایت است که به روش فشرده‌سازی عمیق<sup>۶۷</sup> فشرده شده است. این روش فشرده‌سازی ترکیبی

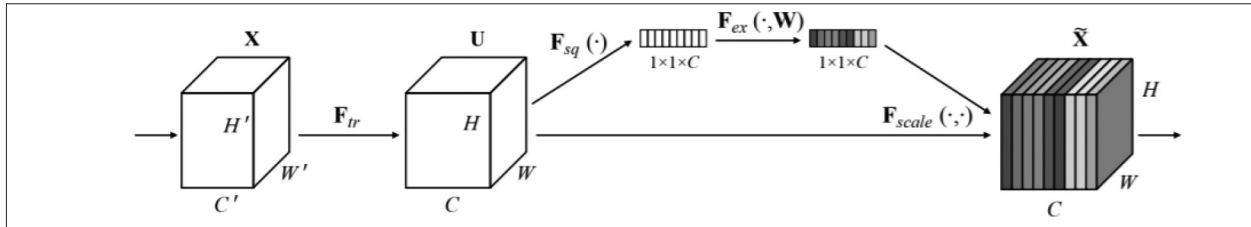
67- Deep Compression

پایاده‌سازی مدل روی FPGAهایی با حجم کوچک. در این مقاله، بلوک‌های سازنده شبکه‌های عصبی پیچشی مانند بلوک رزنت و بلوک اینسپشن را ریزساختار<sup>۶۲</sup> شبکه‌ها، و ساختار کلی آن‌ها را که از چینش ریزساختارها در کنار هم تشکیل می‌شود کلان‌ساختار<sup>۶۳</sup> می‌نامد. ریزساختار ارائه‌شده در این مقاله ماژول آتش<sup>۶۴</sup> نام دارد که با چینش آن‌ها در کنار هم، اسکویزنت ساخته می‌شود.

در این پژوهش چند راهبرد پایه برای طراحی این ماژول جدید ذکر شده است: ۱. جایگزینی پالایه‌های  $3 \times 3$  با پالایه‌های  $1 \times 1$  برای کاهش ۹ برابری پارامترها<sup>۲</sup>. کاهش تعداد لایه‌هایی که ورودی پیش  $3 \times 3$  قرار می‌گیرند، به کمک لایه فشرده‌سازی<sup>۳</sup>. نمونه‌برداری دیرنگام، برای این‌که ورودی لایه‌های پیش اطلاعات بیشتری را شامل باشند که این موضوع در افزایش دقت شبکه موثر است. این قانون که به معنای استفاده کمتر از اندازه گام ۲ در پیش‌ها، و استفاده دیرنگام از لایه ادغام است، موجب می‌شود در بخش عظیم شبکه، نقشه‌های ویژگی اندازه بزرگی داشته باشند. ساختار ماژول آتش (شکل ۲۶) شامل دو بخش فشردن<sup>۶۵</sup> و گستردن<sup>۶۶</sup> است. در بخش فشردن، با استفاده از پیش‌های  $1 \times 1$  تعداد لایه‌های ورودی به بخش گستردن محدود می‌شود. در بخش گستردن، با استفاده از پیش‌های  $1 \times 1$  و  $3 \times 3$  ویژگی‌هایی را استخراج می‌کنند و به‌عنوان خروجی ماژول در کنار هم قرار می‌دهند. برای هم‌اندازه شدن خروجی ماژول، به اندازه یک درایه حاشیه برای ورودی‌های پیش  $3 \times 3$  در نظر گرفته می‌شود.

کلان‌ساختار این شبکه که در شکل ۲۸ قابل مشاهده است، از چینش ماژول‌های آتش و تعداد کمی لایه ادغام شکل گرفته و در انتهای آن ادغام کلی میانگین و softmax برای دسته‌بندی استفاده شده است. پژوهشگران با طراحی و پیاده‌سازی این ساختار، که میزان پارامترهای آن تنها دو درصد الکسنت است، بر روی مجموعه داده ImageNet به

62- Micro Architecture  
63- Macro Architecture  
64- Fire Module  
65- Squeeze  
66- Expand



شکل ۳۱: عملیات فشردن و برانگیختن در ساختار SEnet [۳۶]

است. مقادیر به دست آمده از این تبدیل به عنوان ضریبی در کانال‌های متناظر خود ضرب می‌شوند. شبکه طی مراحل آموزش این ضرایب را به بهترین مقادیر خود نزدیک می‌سازد (شکل ۲۹).

برای سنجش میزان کارایی سازوکار طراحی شده، آن را به نسخه‌های مختلف ساختار رزنت، اینسپشن و موارد دیگر اعمال شده است که در تمامی موارد بهبود دقت را نشان داده است. بهترین نتیجه به دست آمده SEnet-154 است که اعمال سازوکار SE بر نمونه تغییر یافته‌ای از ساختار ResNeXt است.

#### ۴- خلاصه روش‌های بهبود معماری‌ها

##### ۴-۱- افزایش عمق شبکه

افزایش عمق شبکه، یکی از ساده‌ترین ایده‌های مورد استفاده برای افزایش دقت شبکه است که در بسیاری از معماری‌ها از این روش کمک گرفته شده است، خصوصاً در ابتدای کار استفاده از شبکه‌های عصبی پیچشی که توان پردازشی رایانه‌ها رو به افزایش بود روش بسیار موثری به حساب می‌آمد. این روش موجب می‌شود شبکه بتواند ویژگی‌های پیچیده‌تری را یاد بگیرد و در نتیجه دقت افزایش پیدا کند. البته افزایش عمق همواره موجب افزایش دقت نمی‌شود، مسائلی که ممکن است پیش آید یکی افزایش تعداد پارامترهاست که در صورت کافی نبودن حجم دادگان می‌تواند منجر به بیش‌برازش شود. مشکل دیگر این است که وقتی عمق خیلی زیاد شود، حین آموزش شبکه و در مرحله پسرانتشار، گرادینان ناشی از خطاها بسیار کوچک می‌شوند که موجب می‌شود پارامترهای ابتدایی شبکه به‌روزرسانی نشوند و شبکه نتواند به خوبی

از هرس شبکه و چندی سازی<sup>۶۸</sup>، به همراه استفاده از کد هافمن است.

##### ۱۰- SEnet

هدف از ارائه این ساختار، بهبود توان بازنمایی<sup>۶۹</sup> شبکه برای دستیابی به دقت بالاتر است [۳۶]. تمرکز این پژوهش برای این بهبود، تمرکز بر ایجاد ارتباط میان لایه‌های دادگان است. راهکار انتخاب شده اضافه کردن سازوکاری برای مدل کردن ارتباط پویا و خطی مابین لایه‌ها، با استفاده از اطلاعات کلی شبکه است. در این شبکه، تبدیل SE طراحی شده با استفاده از اطلاعات شبکه، یک سازوکار درجه‌ای است که میزان اطلاعات مورد استفاده هر لایه را کنترل می‌کند. این تبدیل قابل اعمال به بلوک‌های سایر معماری‌های شبکه عصبی پیچشی است. تبدیل SE شامل دو بخش فشردن<sup>۷۰</sup> و برانگیختن<sup>۷۱</sup> است. در مرحله فشردن، با استفاده از ادغام کلی میانگین هر کانال از داده تبدیل به یک عدد به عنوان توصیفگر کانال<sup>۷۲</sup> می‌شود. این مقادیر یک شمای کلی از اطلاعات تصویر را می‌سازند. در مرحله برانگیختن، درجه‌بندی تطبیقی<sup>۷۳</sup> مقادیر کانال‌ها با استفاده از توصیفگرهای مرحله قبل انجام می‌شود. درجه‌بندی تطبیقی تعیین وزن‌های کانال‌های داده است که بر اساس ارتباط میان کانال‌ها و میزان اهمیت هر لایه تعیین می‌شود. برای شبیه‌سازی این ارتباطات از دو لایه تمام متصل استفاده می‌شود که ورودی‌های آن توصیفگرهای کانال‌ها هستند و تعداد مقادیر خروجی و ورودی آن برابر است. تابع فعال‌سازی لایه اول ReLU و لایه دوم سیگموئید

- 68- Quantization
- 69- Representational Power
- 70- Squeeze
- 71- Excitation
- 72- Channel Descriptor
- 73- Adaptive Re-Calibration

از عمق افزایش یافته استفاده کند.

#### ۴-۲- ایجاد افزونگی در دادگان

برای بهبود آموزش شبکه، می‌توان تغییراتی بر روی دادگان اولیه انجام داد و دادگان بیشتری را از روی آن‌ها ساخت. این کار علاوه بر افزایش حجم دادگان موجب افزایش تنوع می‌شود. افزایش تنوع دادگان کمک می‌کند شبکه با وابستگی کمتری به حالت ورودی بتواند نتایج را به درستی پیش‌بینی کند.

#### ۴-۳- حذف تصادفی

در این روش در زمان آموزش شبکه، با احتمال مشخص، در هر گذر برخی نورون‌ها را فعال و برخی دیگر را غیرفعال در نظر می‌گیریم، به همین جهت شبکه در هر دور ساختار متفاوتی پیدا می‌کند، و به دلیل عدم اتکای شبکه به تمام ویژگی‌ها برای تصمیم‌گیری، تک تک ویژگی‌های به‌دست آمده قوی‌تر و بهتر انتخاب می‌شوند. ویژگی‌های که می‌توان با تعداد نورون‌های کمتری بر روی آن‌ها تصمیم گرفت، و در واقع ویژگی‌هایی انتخاب می‌شوند که وابستگی کمتری بین نورون‌ها ایجاد می‌کنند. البته این کار می‌تواند در سطح بالاتر از نورون‌ها -مانند بلوک‌های شبکه- هم استفاده شود که در معماری‌های مختلف با آن مواجه خواهیم بود.

#### ۴-۴- تقویت شبکه

منظور از تقویت، این است که چند شبکه آموزش داده شوند و نتیجه نهایی از تصمیم مجموع شبکه‌ها به‌دست آورده شود. ترتیب استفاده از شبکه‌ها و نحوه تعیین نتیجه می‌تواند متفاوت باشد، برای مثال می‌شود چند شبکه را موازی آموزش داد و برای نتیجه از رای‌گیری استفاده نمود.

#### ۴-۵- افزایش عرض شبکه

منظور از افزایش عرض شبکه به‌طور کلی قرار گرفتن تعداد بیشتری از پالایه‌ها در یک لایه از شبکه است، که

موجب می‌شود در هر لایه تعداد ویژگی‌های بیشتری را ذخیره نماییم و افزایش تنوع ویژگی‌ها می‌تواند به افزایش دقت کمک کند. همچنین افزایش عرض پردازش موازی را امکان‌پذیر می‌سازد و موجب می‌شود از توان پردازشی رایانه‌ها بهتر و بهینه‌تر استفاده شود. البته افزایش عرض نیز مانند برخی تکنیک‌های گفته شده می‌تواند روش‌های مختلفی داشته باشد که ساده‌ترین آن افزایش تعداد پالایه‌ها در یک معماری موجود است، و یک مورد پرکاربرد دیگر، افزایش تعداد مسیری‌ها در یک ساختار یا بلوک سازنده ساختار است.

#### ۴-۶- پیچش ۱×۱

این روش که الهام گرفته از معماری شبکه در شبکه [۲۳] است، عملیاتی است که با حفظ طول و عرض نگاشت‌های ویژگی بر روی تعداد آن‌ها موثر است و دادگان را با هم ترکیب می‌کند. این نوع پیچش می‌تواند هم برای افزایش تعداد نگاشت‌ها و هم برای کاهش آن‌ها مورد استفاده قرار گیرد. یکی از کاربردهای آن در جایی است که ما خواهیم خروجی مسیری‌های مختلفی از شبکه را در کنار هم قرار دهیم و برای عملیات مشترکی آماده کنیم که نیاز باشد ابعاد یکسانی داشته باشند، در اینجا می‌توان با استفاده از پیچش ۱×۱ بعد تعداد را با هم برابر کرد. همچنین در برخی ساختارها که نیاز به کاهش ابعاد برای کاهش میزان پردازش هست، از این نوع پیچش استفاده می‌شود زیرا در عین کاهش بعد، بخشی از محتوا بطور کامل از دست نمی‌رود و ترکیبی از آن در لایه‌های دیگر نگاشت موجود است.

#### ۴-۷- ادغام کلی میانگین

این روش برای تبدیل نگاشت‌های ویژگی به یک عدد است که عموماً در لایه‌های آخر شبکه استفاده می‌شود. بدین‌صورت که در لایه‌های آخر، پیش از لایه تصمیم‌گیری، تعداد نقشه‌های ویژگی را به تعداد دسته‌های موجود در دسته‌بندی می‌رسانیم، و سپس میانگین هر کدام



از این نقشه‌های ویژگی را به عنوان مقادیر نهایی به لایه تصمیم‌گیری می‌دهیم تا با روش‌هایی مانند softmax، دسته منتخب را مشخص نماید.

#### ۴-۸- استفاده از دسته‌بند‌های میانی

در مرحله آموزش شبکه، وزن‌های شبکه در مسیر پس‌انتشار خطا به روزرسانی می‌شوند تا به مقادیر بهینه برای تصمیم‌گیری صحیح برسند. با افزایش یافتن عمق شبکه‌ها، خطای انتشار یافته از انتهای شبکه به سمت ابتدای آن مقادیر بسیار کوچکی می‌شود که موجب می‌شود برخی وزن‌ها بخصوص وزن‌های ابتدای شبکه- به روزرسانی نشوند.

یک راه‌حل برای این مشکل استفاده از دسته‌بند‌های میانی است که پیش از رسیدن به انتهای شبکه قرار می‌گیرند و هر کدام بر اساس مقادیری که تا آن بخش از شبکه دیده شده تصمیم‌گیری انجام می‌دهند، در نتیجه می‌توان بر اساس این تصمیم‌ها نیز مقادیر خطایی را محاسبه کرد که در مرحله پس انتشار از این مقادیر برای به‌روزرسانی وزن‌های شبکه استفاده می‌شوند.

#### ۴-۹- شکستن عمل پیچش به پیچش‌های کوچکتر

این روش که در معماری‌های مختلفی استفاده شده است، استفاده از چند پالایه کوچک به صورت متوالی و بدون لایه ادغام در بین آن‌ها، می‌تواند جایگزین پالایه‌هایی با اندازه بزرگتر است. دلایل اصلی برای استفاده از پالایه‌های کوچک -برای مثال  $3 \times 3$ - یکی افزایش دفعات استفاده از غیرخطی‌سازی است که به شبکه کمک می‌کند توابع و ویژگی‌های پیچیده‌تری را یاد بگیرد و دومین دلیل، کاهش تعداد پارامترهاست. برای نمونه در یک لایه پیچش با پالایه‌ای با اندازه  $7 \times 7 \times 3$  که  $C$  تعداد کانال‌ها یا همان نقشه ویژگی‌های ورودی یک لایه است- تعداد پارامترها  $C^2 \times 7^2$  است، اما در شبکه‌ای با سه عمل پیچش  $3 \times 3$  به عنوان جایگزین،  $3(3^2 C^2)$  پارامتر وجود دارد که ۸۱٪ کمتر است. کاهش پارامترها به همگرایی سریع‌تر شبکه و کاهش

بیش‌برازش کمک می‌کند.

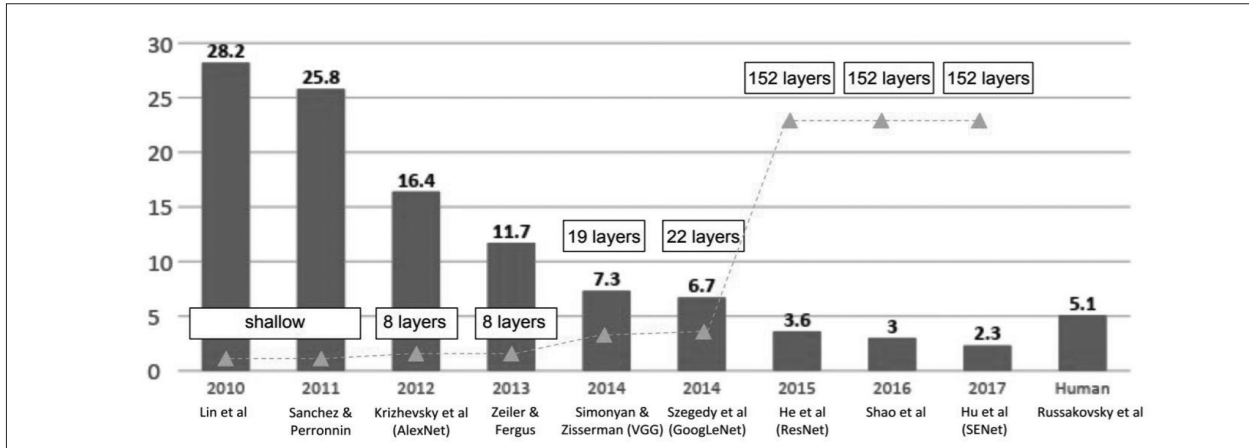
موجب کاهش میزان پردازش در شبکه‌ها و البته موجب افزایش عمق آن‌ها هم می‌شود. این شکست می‌تواند شکست متقارن یا نامتقارن باشد، در بند قبل یک مثال متقارن ذکر شد. در شکست نامتقارن می‌شود پیچش  $n \times n$  را به دو پیچش  $n \times 1$  و  $1 \times n$  شکست. برای مثال به جای یک فیلتر  $3 \times 3$  با ۹ پارامتر، از دو پیچش  $1 \times 3$  و  $3 \times 1$  -هر کدام شامل ۳ پارامتر- استفاده می‌شود که با ۳۳ درصد کاهش تعداد پارامترها همراه است. البته در مقاله [۲۶] ذکر شده است که این کار در لایه‌های ابتدایی شبکه عملکرد خوبی را نشان نمی‌دهد و در لایه‌هایی با نگاشت‌های ویژگی به ابعاد ۱۲ الی ۲۰ نتیجه خوبی را به همراه دارد.

#### ۴-۱۰- استفاده از اتصال‌های میانبر

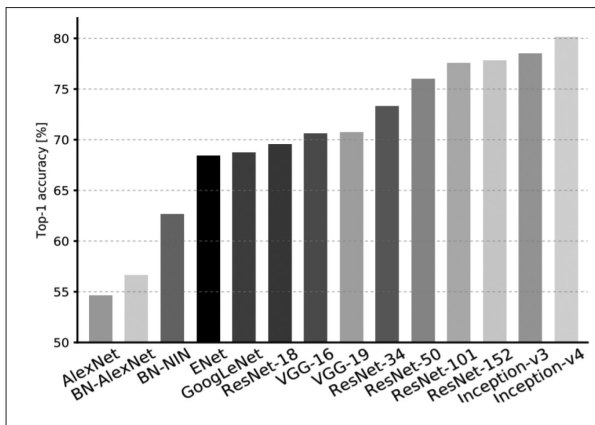
یکی دیگر از راه‌حل‌های ارائه شده برای رفع مشکل محو گرادیان در شبکه‌های عمیق، استفاده از اتصال‌های میانبر است که ساختارهای مبتنی بر باقیمانده را ایجاد می‌کنند. اگر خروجی شبکه در حالت معمول، تابع  $H(x)$  باشد، ما با افزودن میان‌بر به بلوک سازنده معماری، مقدار  $x$  یا همان ورودی بلوک را به خروجی آن اضافه می‌کنیم، در نتیجه شبکه باید بتواند به کمک پارامترهای خود، تابع  $H(x)-x$  را بسازد که همان بازنمایی باقیمانده است.

#### ۵- مقایسه و جمع‌بندی

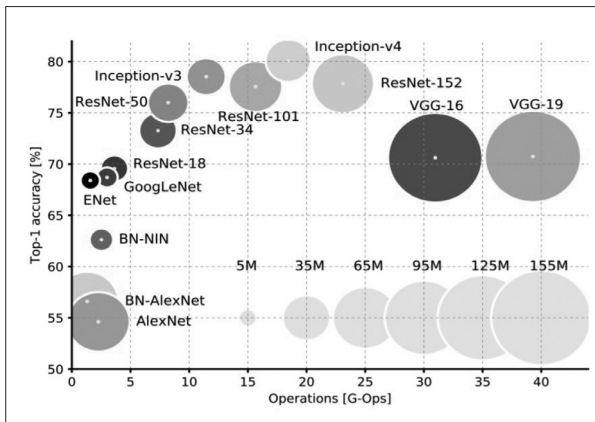
در این پژوهش مروری داشتیم بر ساختار شبکه عصبی پیچشی و لایه‌های اصلی تشکیل‌دهنده آن، سپس معماری‌های مهم شبکه عصبی پیچشی که تغییرات جدی‌تری در این زمینه ایجاد کرده بودند، بررسی و نکات اصلی تفاوت آن‌ها ذکر شد. تلاش‌های انجام شده برای افزایش دقت عمدتاً در راستای افزایش عمق، تغییر ساختار در جهت بهبود پس‌انتشار خطا و ایجاد انعطاف در شبکه بوده است. برخی پژوهشگران نیز با هدف کاهش تعداد پارامترها تلاش در تغییر ساختار برای بازاستفاده از ویژگی‌ها دارند. این حوزه، یکی از حوزه‌ها فعال پژوهشی



شکل ۳۲: خطای تشخیص و تعداد لایه‌های معماری‌های برنده چالش دسته‌بندی ImageNet [۳۷]



شکل ۳۳: مقایسه دقت اولین تشخیص صحیح معماری‌های مختلف بر روی مجموعه داده ImageNet [۳۸]



شکل ۳۴: مقایسه حجم پردازش و دقت اولین تشخیص صحیح معماری‌های مختلف بر روی مجموعه داده ImageNet [۳۸]

ImageNet که در شکل ۳۱ نتایج آن آمده است و دومی میزان پردازش موردنیاز برای آموزش هر یک از ساختارها در کنار دقت تشخیص است (شکل ۳۲). به‌عنوان جمع‌بندی معماری‌های مختلف، خلاصه‌ای از

است که هر روزه نوآوری‌های تازه‌ای را در آن شاهد هستیم. ارائه ساختارهای جدید، بهبود ساختارهای موجود، کاهش تعداد پارامترها و موارد دیگر می‌توانند حوزه‌های فعالیت برای پژوهشگران علاقه‌مند باشد.

به‌عنوان بخشی از جمع‌بندی این مقاله، مقایسه‌ای بین معماری‌های معرفی‌شده ارائه می‌شود که برای این هدف، می‌توان به میزان کارایی این معماری‌ها برای چالش دسته‌بندی ImageNet [۱] اشاره نمود که تا سال ۲۰۱۷ جریان داشته است و هر ساله برای حل این چالش معماری‌های جدید شکل گرفته و معرفی شده‌اند. همان‌طور که در شکل ۳۲ قابل مشاهده است، برنده این چالش از سال ۲۰۱۲ و پس از آن به ترتیب الکس‌نت، VGGnet، ZFnet و گوگل‌نت (به ترتیب در چالش دسته‌بندی و دیگری در چالش تشخیص محل شیء ۷۴) و رزنت بوده است، برنده سال ۲۰۱۶ نیز با چینی‌ها چند نوع از معماری‌های پیشین در کنار هم و جمع نتایج به این میزان خطا دست یافته است. خطای گزارش‌شده در شکل ۳۰ مربوط به چالش تشخیص دسته صحیح از میان ۵ حدس اول است. این چالش در سال ۲۰۱۷ با ارائه ساختار SENet و دستیابی به خطای ۲٫۳ درصد به پایان رسیده و دیگر ادامه پیدا نکرده است.

همچنین برای مقایسه می‌توان به پژوهش [۳۸] اشاره نمود که معماری‌ها را از دو منظر بررسی کرده است. یکی دقت اولین تشخیص صحیح، بر روی مجموعه داده

جدول ۲: مقایسه معماری‌های مختلف شبکه‌های عصبی پیچشی و نوآوری‌های آن‌ها

سال انتشار مقالات	ساختار برتر دسته			خلاصه نوآوری‌ها نام ساختار	نام معماری
	خطا روی ImageNet	تعداد پارامترها	عمق شبکه		
۱۹۹۸-۱۹۹۰	-	۶۰۰۰۰	۷	لینت-۵	لینت [۱۱] [۱۷]
۲۰۱۲	16.4% *	۶۲ میلیون	۸	الکسنت	الکسنت [۲۰]
۲۰۱۴	16.1% *	ذکر نشده	۸	ZFnet با ۵۱۲، ۵۱۲ و ۱۰۲۴ فیلتر در لایه ۳ و ۴ و ۵	ZFnet [۲۲]
۲۰۱۳	-	ذکر نشده	۴	شبکه در شبکه	شبکه در شبکه [۲۳]
۲۰۱۴	7.1%	۱۴۴ میلیون	۱۹	VGG-E	VGGnet [۲۶]
۲۰۱۶-۲۰۱۴	3.7%	۵۵ میلیون [۴۰]	۱۶۴ [۳۹]	اینسپین-رزنت ۲	گوگل نت (اینسپین) [۲۴] [۲۷] [۲۸]
۲۰۱۵	3.6%	۶۰ میلیون [۴۱]	۱۵۲	رزنت ۱۵۲	رزنت [۲۵]
۲۰۱۵	-	۲.۳ میلیون	۱۹	هایوی ۲	هایوی نت [۲۹]
۲۰۱۶	4.8%	۶۴.۲ میلیون [۳۱]	۲۰۰	رزنت ۲۰۰ پیش‌فعال شده ۱۵	پیش‌فعال سازی رزنت [۳۰]
۲۰۱۶	4.9%	ذکر نشده	۱۲۰۲	ندارد	شبکه با عمق تصادفی [۳۲]
۲۰۱۶	6.03%	۶۸.۹ میلیون	۵۰	WRN-۵۰-۲	رزنت عریض [۳۱]
۲۰۱۷	4.4%	۸۳.۶ میلیون	۱۰۱	ResNeXt ۴×۶۴-۱۰۱	ResNeXt [۳۳]
۲۰۱۷	5.29%	۴-۳۳ میلیون	۲۶۴	دنسنت-۲۶۴ (نرخ رشد ۳۲)	دنسنت [۳۴]
۲۰۱۶	17.5%	۱.۲ میلیون [۴۱]	۶۹ [۴۱]	اسکوویزنت + اتصال ساده	اسکوویزنت [۳۵]
۲۰۱۷	3.58% *	۱۴۶ میلیون [۴۰]	۱۵۴	SEnet-154	SEnet [۳۶]

ایده‌های بهبود رزنت

1- ResNet200 (pre-act)

529, no. 7587, pp. 484–489, 2016.

- [11] LeCun, Yann, Boser, B. E., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W. E., Jackel, L. D., "Handwritten Digit Recognition with a Back-Propagation Network," in *Advances in Neural Information Processing Systems 2*, D. S. Touretzky, Ed. Morgan-Kaufmann, 1990, pp. 396–404.
- [12] HUBEL, D. H., WIESEL, T. N., "Receptive fields of single neurones in the cat's striate cortex.," *J. Physiol.*, vol. 148, no. 3, pp. 574–91, Oct. 1959.
- [13] "Plato's Cave: Hubel & Wiesel." [Online]. Available: <http://cns-alumni.bu.edu/~slehar/webstuff/pcave/hubel.html>. [Accessed: 12-Sep-2019].
- [14] "CS231n Convolutional Neural Networks for Visual Recognition." [Online]. Available: <http://cs231n.github.io/convolutional-networks/>. [Accessed: 12-Sep-2019].
- [15] Lecun, Y., Bengio, Y., Hinton, G., "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [16] "A Beginner's Guide To Understanding Convolutional Neural Networks – Adit Deshpande – Engineering at Forward | UCLA CS '19." [Online]. Available: <https://adeshpande3.github.io/adeshpande3.github.io/A-Beginner's-Guide-To-Understanding-Convolutional-Neural-Networks/>. [Accessed: 12-Sep-2019].
- [17] LeCun, Y., LeCun, Y., Jackel, L., Bottou, L., Brunot, A., Cortes, C., Denker, J., Drucker, H., Guyon, I., Müller, U., Säcker, E., Simard, P., Vapnik, V., "Comparison of Learning Algorithms for Handwritten Digit Recognition," *Int. Conf. Artif. NEURAL NETWORKS*, pp. 53–60, 1995.
- [18] Lecun, Y., "Gradient-Based Learning Applied to Document Recognition," *proc. IEEE*, 1998.
- [19] "The EMNIST Dataset | NIST." [Online]. Available: <https://www.nist.gov/node/1298471/emnist-dataset>. [Accessed: 09-Aug-2019].
- [20] Krizhevsky, A., Sutskever, I., Hinton, G. E., "Imagenet classification with deep convolutional neural networks," *Adv. Neural Inf. Process. Syst.*, p. 2012.
- [21] Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., Salakhutdinov, R. R., "Improving neural networks by preventing co-adaptation of feature detectors," Jul. 2012.
- [22] Zeiler, M. D., Fergus, R., "Visualizing and Understanding Convolutional Networks," Springer, Cham, 2014, pp. 818–833.
- [23] Lin, M., Chen, Q., Yan, S., "Network In Network," *CoRR*, vol. abs/1312.4, pp. 1–10, 2013.
- [24] Simonyan, K., Zisserman, A., "Very Deep Convolutional Networks for Large-Scale Image Recognition," pp. 1–14, 2014.
- [25] He, K., Zhang, X., Ren, S., Sun, J., "Deep Residual Learning for Image Recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [26] Simonyan, K., Zisserman, A., "Very Deep Convolutional Networks for Large-Scale Image Recognition," pp. 1–14, 2014.
- [27] Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna,

مشخصات شبکه‌ها و معماری برتر ذکر شده در مقالات آن دسته در جدول ۲ آورده شده است. در این جدول، در بخش ساختار برتر، صرفاً ساختارهای تکی در نظر گرفته شده است و مواردی که با گروه کردن چند ساختار به خطای کمتری می‌رسند در جدول نیامده است. خطای اعلام شده در این جدول، مربوط به ۵ تشخیص است. همچنین مواردی که با \* مشخص شده‌اند، بر روی مجموعه داده آزمایشی گزارش شده‌اند و در سایر موارد خطای گزارش شده بر روی داده اعتبارسنجی<sup>۵</sup> است.

## مراجع

- [1] "ImageNet." [Online]. Available: <http://www.image-net.org/>. [Accessed: 03-Aug-2019].
- [2] Wang M., Deng W., "Deep Face Recognition: {A} Survey," *CoRR*, vol. abs/1804.0, 2018.
- [3] Wang T., Wu D. J., Coates A., Ng A. Y., "End-to-end text recognition with convolutional neural networks," in *Proceedings of the 21st International Conference on Pattern Recognition (ICPR2012)*, 2012, pp. 3304–3308.
- [4] Yamashita R., Nishio M., Do R. K. G., Togashi K., "Convolutional neural networks: an overview and application in radiology," *Insights Imaging*, vol. 9, no. 4, pp. 611–629, 2018.
- [5] Ekici S., Jawzal H., "Breast cancer diagnosis using thermography and convolutional neural networks," *Med. Hypotheses*, vol. 137, p. 109542, 2020.
- [6] Schwendicke F., Golla T., Dreher M., Krois J., "Convolutional neural networks for dental image diagnostics: A scoping review," *J. Dent.*, vol. 91, p. 103226, 2019.
- [7] Gruszczyński W., Puniach E., Ćwiakła P., Matwij W., "Application of convolutional neural networks for low vegetation filtering from data acquired by UAVs," *ISPRS J. Photogramm. Remote Sens.*, vol. 158, pp. 1–10, 2019.
- [8] Brahimi M., Arsenovic M., Laraba S., Sladojevic S., Boukhalfa K., Moussaoui A., "Deep Learning for Plant Diseases: Detection and Saliency Map Visualisation," in *Human and Machine Learning: Visible, Explainable, Trustworthy and Transparent*, Zhou J., Chen F., Eds. Cham: Springer International Publishing, 2018, pp. 93–117.
- [9] Brodrick P. G., Davies A. B., Asner G. P., "Uncovering Ecological Patterns with Convolutional Neural Networks," *Trends Ecol. Evol.*, vol. 34, no. 8, pp. 734–745, 2019.
- [10] Silver D., Huang A., Maddison Ch. J., Guez A., Sifre L., van den Driessche G., Schrittwieser J., Antonoglou I., Panneershelvam V., Lanctot M., Dieleman S., Grewe D., Nham J., Kalchbrenner N., Sutskever I., Lillicrap T., Leach M., Kavukcuoglu K., Graepel Th., Hassabis D., "Mastering the game of Go with deep neural networks and tree search," *Nature*, vol.

tion (CVPR), 2017, pp. 2261–2269.

[35] Iandola, F. N., Han, S., Moskewicz, M. W., Ashraf, K., Dally, W. J., Keutzer, K., Han, S., Dally, W. J., “SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <1MB model size,” CoRR, vol. abs/1602.0, pp. 1–13, 2016.

[36] Hu, J., Shen, L., Albanie, S., Sun, G., Wu, E., “Squeeze-and-Excitation Networks,” CoRR, vol. abs/1709.0, pp. 7132–7141, Sep. 2017.

[37] Li, F.-F., Johnson, J., Yeung, S., “Lecture 9: CNN Architectures,” 2019. [Online]. Available: [http://cs231n.stanford.edu/slides/2019/cs231n\\_2019\\_lecture09.pdf](http://cs231n.stanford.edu/slides/2019/cs231n_2019_lecture09.pdf). [Accessed: 03-Aug-2019].

[38] Canziani, A., Paszke, A., Culurciello, E., “An Analysis of Deep Neural Network Models for Practical Applications,” CoRR, vol. abs/1605.0, 2016.

[39] “Pretrained Inception-ResNet-v2 convolutional neural network - MATLAB inceptionresnetv2.” [Online]. Available: <https://www.mathworks.com/help/deeplearning/ref/inceptionresnetv2.html>. [Accessed: 03-Jan-2020].

[40] “State-of-the-art table for Image Classification on ImageNet.” [Online]. Available: <https://paperswithcode.com/sota/image-classification-on-imagenet>. [Accessed: 03-Jan-2020].

[41] “Image - Wolfram Neural Net Repository.” [Online]. Available: <https://resources.wolframcloud.com/NeuralNetRepository/inputdomain/Image/>. [Accessed: 03-Jan-2020].

Z., “Rethinking the Inception Architecture for Computer Vision,” in 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 2818–2826.

[28] Szegedy, C., Ioffe, S., Vanhoucke, V., Alemi, A. A., “Inception-v4, inception-ResNet and the impact of residual connections on learning,” Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence. AAAI Press, pp. 4278–4284, 2017.

[29] Srivastava, R. K., Greff, K., Schmidhuber, J., “Highway Networks,” CoRR, vol. abs/1505.0, May 2015.

[30] He, K., Zhang, X., Ren, S., Sun, J., “Identity Mappings in Deep Residual Networks,” CoRR, vol. abs/1603.0, pp. 630–645, 2016.

[31] Zagoruyko, S., Komodakis, N., “Wide Residual Networks,” CoRR, vol. abs/1605.0, 2016.

[32] Huang, G., Sun, Y., Liu, Z., Sedra, D., Weinberger, K. Q., “Deep Networks with Stochastic Depth,” in Computer Vision -- ECCV 2016, 2016, pp. 646–661.

[33] Xie, S., Girshick, R., Dollár, P., Tu, Z., He, K., “Aggregated Residual Transformations for Deep Neural Networks,” in 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, pp. 5987–5995.

[34] Huang, G., Liu, Z., Maaten, L. van der, Weinberger, K. Q., “Densely Connected Convolutional Networks,” in 2017 IEEE Conference on Computer Vision and Pattern Recogni-

## تراوش های ذهنی

۲۵ شیوه نگرش به هوش مصنوعی

ترجمه: ابراهیم نقیبزاده مشایخ

ست لوید

جودیتا بزل

استوارت راسل

جرج دایسون

دانیل کینت

رادنی بروکز

فرانک ویلیک

مکس تگمارک

ویراستار:

جان

بروکمن

ونکی زاماکیشیان

الکس «سندی» بنتاند

هانس اولریش اوربست

آلیسون کوپنیک

بیتر کالیسون

جرج مکدونالد جرج

کارولینا جونز

استفن ولفرام

جان تالین

استیون بینگر

دیوید تویچ

تام گریفیث

آکا دراگان

کریس اندرسون

دیوید کایزر

نیل کوشنفلد

دانیل هیلیس

ISI

انجمن انفورماتیک ایران

## جدیدترین کتاب

از انتشارات انجمن انفورماتیک ایران

منتشر شد!

# تراوش های ذهنی

تهیه کتاب از دفتر انجمن انفورماتیک ایران

(۶۶۴۱۲۸۶۱) و فروشگاه اینترنتی چاره

[www.chare.ir](http://www.chare.ir)

قیمت ۴۰/۰۰۰ تومان