

# بازیابی ریزدانه‌ای تصویر به کمک شبکه خودتوجهی مکانی و مکانیسم برجسته‌سازی

سید نیما سید آقا یزدی<sup>۱</sup>، نام و نام خانوادگی نویسنده دوم<sup>۲</sup>، نام و نام خانوادگی نویسنده سوم<sup>۳</sup>

<sup>۱</sup> دانشجوی کارشناسی ارشد هوش مصنوعی، دانشکده فنی مهندسی، دانشگاه آزاد تهران جنوب، تهران،  
nima.yazdi۹۳@gmail.com

<sup>۲</sup> رتبه علمی نویسنده در صورت تمایل، گروه آموزشی یا واحد سازمانی مربوطه، نام سازمان، شهر  
آدرس پست الکترونیکی

<sup>۳</sup> رتبه علمی نویسنده در صورت تمایل، گروه آموزشی یا واحد سازمانی مربوطه، نام سازمان، شهر  
آدرس پست الکترونیکی

## چکیده

بازیابی تصاویر، دسته‌بندی دقیق تصاویر، با استفاده از شباهت‌ها و تفاوت‌های موجود در بافت، رنگ، فرم و سایر ویژگی‌های تصویر است. بازیابی تصویر برای پرس‌وجوی مبتنی بر تصویر شامل رویکردهای متفاوتی است که می‌توان آن‌ها را در سه دسته عمده بیان نمود: بازیابی تصویر مبتنی بر طرح، بازیابی تصویر مبتنی بر محتوا و بازیابی تصویر مبتنی بر ریزدانه. در این مقاله شبکه خودتوجهی مکانی پیشنهاد شده است که شامل دو جزء اصلی می‌باشد. ابتدا یک شبکه عصبی کانولوشنی به‌عنوان استخراج‌کننده ویژگی پیاده‌سازی می‌شود که ویژگی‌های اولیه را از تصاویر ورودی از طریق چندین لایه کانولوشن استخراج می‌کند. سپس ماژول خودتوجهی مکانی با استفاده از مکانیسم توجه، ویژگی‌های جدید را ذخیره می‌کند. یکی از مشکلات روش استفاده از شبکه خودتوجهی مکانی آن است که تصویر ورودی، با ویژگی‌های با اهمیت کم‌تر بررسی می‌شود و ممکن است بخش‌های حاشیه‌ای در نتیجه نهایی عملکرد شبکه، تأثیرگذار باشند. در این مقاله، روش XRAI برجسته‌سازی پیشنهاد شده است. این روش با امتیاز ۸۸ درصد توانسته است نتیجه‌ای قابل توجه نسبت به سایر روش‌های بازیابی تصویر داشته باشد.

## کلمات کلیدی

بازیابی تصویر - بازیابی ریزدانه‌ای تصویر - بازیابی تصویر مبتنی بر محتوا - شبکه خودتوجهی مکانی - مکانیسم برجسته‌سازی

## ۱- مقدمه

پیش می‌آید. از جمله آنکه ویژگی‌های استخراج شده با ادراک انسان فاصله معنایی بسیاری داشتند. اما با انتخاب و استخراج درست ویژگی‌های مورد محاسبه، این فاصله کمتر به چشم آمده است. به‌گونه‌ای که اکنون با نیاز به بررسی دقیق‌تر دسته‌بندی تصاویر، بازیابی تصاویر ریزدانه‌ای<sup>۲</sup> معرفی شده است که در پیدا کردن ویژگی‌های مشابه، تا حد ادراک انسان رفتار می‌کند.

بازیابی تصویر برای پرس‌وجوی مبتنی بر تصویر، شامل رویکردهای متفاوتی است که می‌توان آن‌ها را در سه دسته عمده بیان نمود: بازیابی تصویر مبتنی بر طرح، بازیابی تصویر مبتنی بر محتوا و بازیابی تصویر مبتنی بر ریزدانه. این دسته بندی‌ها هر کدام دارای زیرروش‌های مختلفی هستند که می‌توان آن‌ها را با توجه به انواع

بازیابی تصاویر، دسته‌بندی دقیق تصاویر، با استفاده از شباهت‌ها و تفاوت‌های موجود در بافت، رنگ، فرم و سایر ویژگی‌های تصویر است. این شاخه از علم پردازش تصویر، برای اولین بار در سال ۱۹۷۰ با رویکرد مبتنی بر متن<sup>۱</sup> معرفی گردید. پس از آن رویکردی متفاوت، با عنوان مبتنی بر محتوا<sup>۲</sup> معرفی شد که بر اساس ویژگی‌های استخراج شده از تصاویر، کار می‌کرد. این رویکرد به‌سرعت جایگزین رویکرد پیشین شد و در حوزه‌های پزشکی، احراز هویت، پیشگیری از وقوع جرم، امنیت محیط و... مورد استفاده قرار گرفت. در این میان چالش‌های بسیاری به هنگام استفاده از روش‌های مبتنی بر این رویکرد،

سیامی کانولوشنی پیشنهاد می‌شود. ابتدا، تکه‌های ضایعه برای ایجاد دو مجموعه‌های داده، برش داده می‌شوند و جفت‌های دوتکه دلخواه یک مجموعه‌داده پیچ-جفت را تشکیل می‌دهند. دوم، این مجموعه‌داده پیچ-جفت برای آموزش یک شبکه استفاده می‌شود. سوم، یک پیچ آزمایشی به‌عنوان یک پرس‌وجو در نظر گرفته می‌شود. فاصله بین این پرس‌وجو و بیست مورد در هر دو مجموعه‌داده با استفاده از شبکه عصبی کانولوشنی سیامی آموزش دیده محاسبه می‌شود. موارد نزدیک به پرس‌وجو برای ارائه پیش‌بینی نهایی با بالاترین امتیاز استفاده می‌شود. در [۹] روشی پیشنهاد می‌شود که از قدرت شبکه‌های عصبی کانولوشن برای پیش‌بینی دسته‌بندی تصویر ورودی برای همه کلاس‌های خروجی و بازیابی تصاویر با استفاده از تابع فاصله تغییر یافته در فضای ویژگی موجه، استفاده می‌کند. در [۱۰] بافت کانال و اطلاعات توالی مکانی برای بازیابی مبتنی بر محتوا مورد تمرکز قرار می‌گیرند. ابتدا یک مدل عمیق پیشنهاد می‌شود که هدف آن استنباط نقشه‌های توجه، در امتداد ابعاد کانال و مکان است. با بهبود ماژول‌های توجه کانال و توجه مکانی و کاوش ترانسفورماتور<sup>۱۲</sup>، توانایی ساخت و درک مدل افزایش می‌یابد. در [۱۱] با اشاره به روش‌هایی که با خطای ویژگی‌های عمومی به استخراج ویژگی‌های متمایزتر کمک می‌کنند، یک تابع محاسبه خطای جدید به نام خطای متمرکز سخت ارائه می‌شود. این تابع در استخراج ویژگی برای تمایز در تقسیم مشابه‌ترین دسته‌ها کمک می‌کند. در [۱۲] یک شبکه ترکیبی مبتنی بر خودتوجهی<sup>۱۳</sup> برای یادگیری بازنمایی‌های رایج داده‌های رسانه‌های مختلف<sup>۱۴</sup> پیشنهاد می‌شود. به طور خاص، ابتدا از یک لایه خودتوجهی محلی برای یادگیری فضای توجه مشترک بین داده‌های رسانه‌های مختلف استفاده می‌شود. سپس یک روش الحاق شباهت برای درک رابطه محتوایی بین ویژگی‌ها پیشنهاد می‌شود. برای بهبود بیشتر استحکام مدل، یک کدگذاری موقعیت محلی را یاد می‌گیرد تا روابط مکانی بین ویژگی‌ها را ثبت کند؛ در [۱۳] یک چارچوب سبک‌تر برای نمونه‌برداری تدریجی از قطعات متمایز، جهت یادگیری جزئیات ارائه می‌شود. در این روش ابتدا شیء از تصویر اصلی تقویت شده و سپس یک نمونه‌برداری خودتطبیقی برای شناسایی بیشتر منطقه تقویت شده انجام می‌گردد. پس این چارچوب می‌تواند از کل به شیء و از شیء به جزئیات برسد.

در دسته روش‌های مبتنی بر ریزدانه، در [۱۴] انتخاب توصیف‌گرهای عمیق مفید به خوبی به تشخیص تصویر با دانه‌ریز کمک می‌کند. به طور خاص، یک مدل جدید شبکه عصبی کانولوشنی ماسک دار<sup>۱۵</sup>، بدون لایه‌های کاملاً متصل پیشنهاد شده است. بر اساس حاشیه‌نویسی‌های بخش، مدل پیشنهادی شامل یک شبکه کاملاً کانولوشنی برای مکان‌یابی قسمت‌های متمایز و مهم‌تر از آن تولید ماسک‌های جسم/قطعه وزن‌دار برای انتخاب توصیف‌گرهای کانولوشنی مفید و معنادار است. پس از آن، یک مدل سه‌جریانی برای جمع‌آوری توصیف‌گرهای انتخاب شده در سطح جسم و قطعه به طور هم‌زمان ساخته می‌شود. با کنار گذاشتن پارامتر لایه‌های کاملاً متصل اضافی، این شبکه دارای ابعاد کوچک و سرعت استنتاج کارآمد در مقایسه با سایر روش‌های ریزدانه است. در [۱۵] یک روش تخمین ریزدانه، برای تخمین نمره زیبایی‌شناسی<sup>۱۶</sup> پیشنهاد می‌شود و مکانیسم‌های توجه موقعیت و کانال را برای افزایش ترکیب ویژگی‌های زیبایی‌شناسی

استخراج ویژگی، پردازش ویژگی و طبقه‌بندی تصاویر، دسته‌بندی کرد. در ادامه به بررسی نمونه‌هایی از این دسته‌ها پرداخته شده است.

در دسته روش‌های مبتنی بر طرح، در [۱] یک مدل بازیابی تصویر مبتنی بر طرح با نام صفرشات<sup>۱۷</sup> را بررسی می‌کند که در آن دسته‌های آزمایشی در مرحله آموزش ظاهر نمی‌شوند. سعی می‌شود از طریق گسستگی نامتقارن<sup>۱۸</sup> به بازیابی آگاهانه از ساختار رسید. که در آن ویژگی‌های تصویر به ویژگی‌های ساختار و ویژگی‌های ظاهری تفکیک می‌شوند درحالی‌که ویژگی‌های طرح تنها به فضای ساختار، پیش‌بینی می‌شوند. از طریق جداسازی ساختار و فضای ظاهری، ترجمه دامنه دوجهته بین حوزه طرح و حوزه تصویر انجام می‌شود. در [۲] بازیابی تصویر مبتنی بر طرح به‌عنوان یک فرایند درشت به ریز فرموله شده است و یک مدل رتبه‌بندی متقابل با نام آبشاری عمیق<sup>۱۹</sup> پیشنهاد می‌شود که می‌تواند از تمام اطلاعات چندوجهی مفید در طرح‌ها و تصاویر حاشیه‌نویسی بهره‌برداری کند و کارایی بازیابی را بهبود بخشد. هدف ساختن بازنمایی‌های عمیق برای طرح‌ها، تصاویر، توضیحات و یادگیری همبستگی‌های عمیق بهینه شده در چنین حوزه‌های مختلف، متمرکز است؛ بنابراین برای یک طرح ورودی داده شده، تصاویر مربوط به آن با شباهت‌های ریز در سطح نمونه، در یک دسته خاص می‌توانند برگردانده شوند.

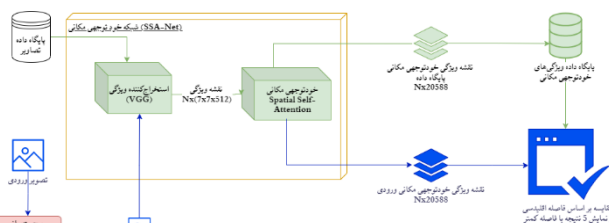
در دسته روش‌های مبتنی بر محتوا، در [۳] استفاده از جنگل‌های مسیر بهینه (بدون نظارت و با نظارت) و رویکردهای یادگیری فعال را برای بازخورد مرتبط در سیستم‌های بازیابی تصویر پزشکی مبتنی بر محتوا بررسی می‌کند. آموزنده‌ترین تصاویری که با رویکرد یادگیری فعال انتخاب می‌شوند، آن‌هایی هستند که بهترین تعادل را بین شباهت (با تصویر پرس‌وجو) و درجات خاصی از تنوع و عدم قطعیت ارائه می‌دهند. در [۴] هدف ایجاد یک روش بازیابی تصویر برای مشخص کردن دسته‌بندی یک محصول است و یک مدل شبکه کانولوشنی سیامی<sup>۲۰</sup> پیشنهاد می‌شود که شامل برجسب‌های دسته و مورد در آموزش برای تولید ویژگی آگاه از دسته است. این مدل با اصلاح رویه آموزشی همراه می‌باشد که به طور هم‌زمان دسته و برجسب مورد را یاد می‌گیرد.

در [۵] یک روش بازیابی متقابل رسانه‌ای مبتنی بر ترکیب چند ویژگی<sup>۲۱</sup> پیشنهاد می‌شود. این روش قادر به ادغام چندین ویژگی برای ارتقای درک معنایی و اتخاذ یادگیری متخاصم<sup>۲۲</sup> برای بهبود بیشتر دقت بازنمایی زیرفضای عمومی است. سپس از شباهت در همان فضا برای مرتب‌سازی نتایج بازیابی استفاده می‌شود. در [۶] یک روش بازیابی تصویر مبتنی بر محتوا پیشنهاد می‌شود. در مرحله توصیف تصویر، این روش ابتدا توصیفگر ریزساختار سنتی را اصلاح می‌کند تا رابطه مستقیم بین ویژگی‌های شکل و بافت و بین ویژگی‌های رنگ و بافت را به تصویر بکشد. سپس هیستوگرام الگوهای باینری محلی یکنواخت<sup>۲۳</sup> تصویر را استخراج می‌کند تا اطلاعات تفاوت رنگ را به تصویر بکشد. در [۷] یک چارچوب چند وظیفه‌ای مبتنی بر جداسازی و بازسازی ویژگی<sup>۲۴</sup> برای بازیابی متقابل وجهی بر اساس روش‌های رایج یادگیری مکانی پیشنهاد می‌شود که ماژول جداسازی ویژگی را برای مقابله با عدم تقارن اطلاعات بین روش‌های مختلف معرفی می‌کند. در [۸] بازیابی تصویر مبتنی بر محتوا با یک شبکه عصبی

ترکیب می‌کند. با آموزش شبکه رگرسیون. جدا از شبکه طبقه‌بندی، طبقه‌بندی وظیفه رگرسیون را تکمیل می‌کند. محققان به استفاده از میانگین مربع خطا<sup>۱۷</sup> به عنوان معیار ارزیابی اصلی عادت کرده‌اند، که در اندازه‌گیری خطای هر بازه ناکافی است. به منظور در نظر گرفتن کامل تصاویر، بخش‌های مختلف امتیاز زیبایی‌شناسی، به جای تمرکز بر بخش‌های این نمره‌ها به دلیل عدم تعادل مجموعه داده‌ها، یک معیار ارزیابی جدید به نام خطاهای میانگین مربع تقسیم شده<sup>۱۸</sup> برای اثبات مزایا پیشنهاد می‌شود. در [۱۶] که روی طبقه‌بندی تصاویر تمرکز دارد، یک سیستم یادگیری نیمه نظارت تهیه شده است. در این روش یک مکانیسم توجه تعاملی ریزدانه‌ای تعبیه شده که در ابتدا از تصاویر برجسب‌دار استفاده کرده و به تهیه بردارهای احتمالی حاصل از این تصویر، می‌پردازد. سپس داده‌های آموزشی بدون برجسب را با این بردارها مقایسه و طبقه‌بندی می‌کند. در [۱۷] روش یادگیری هش<sup>۱۹</sup> با دو مشکل بررسی می‌شود: ۱- ویژگی‌های با ابعاد کم، فرایند بازیابی را تسریع می‌بخشد اما به دلیل از دست رفتن اطلاعات، دقت را کاهش می‌دهند. ۲- تصاویر ریزدانه منجر به ایجاد کدهای هش جستجوی یکسان در خوشه‌های مختلف در فضای پنهان پایگاه داده می‌شوند. از این رو به یک شبکه پاک‌کننده توجه مبتنی بر ثبات ویژگی<sup>۲۰</sup> پرداخته می‌شود. برای مشکل نخست، از یک ماژول پاک‌کردن ناحیه انتخاب شده<sup>۲۱</sup> استفاده می‌کند که با پوشش تطبیقی برخی از مناطق تصاویر خام، شبکه را در برابر تفاوت‌های ظریف ریزدانه‌ای مقاوم می‌کند. بدین ترتیب کدهای هش متمایزتری در پایگاه داده هش ذخیره می‌شوند. سپس برای پایدارتر کردن رابطه بین کد هش تصویر ورودی و کد هش پایگاه داده، از ماژول افزایش خطای رابطه مکانی<sup>۲۲</sup> استفاده می‌کند. در [۱۸] یک شبکه بازیابی و استخراج اطلاعات متمایز به نام DRE-Net پیشنهاد می‌شود که با مشکل تشخیص تصویر با رزولوشن پایین رسیدگی می‌کند. این شبکه از دو شبکه فرعی تشکیل شده است: ۱- زیر شبکه بازیابی اطلاعات متمایز ریز ۲۳- ۲- زیر شبکه شناسایی با رابطه معنایی خطای تقطیر<sup>۲۴</sup>. ماژول اول با استفاده از ویژگی‌ها، به بازیابی جزئیات بافت حیاتی پیکسل‌ها کمک می‌کند. ماژول دوم به روابط صحیح بین هر دو پیکسل در نقشه ویژگی می‌پردازد. پس ماژول دوم می‌تواند به ماژول اول برای پیدا کردن جزئیات دقیق و قابل اعتماد کمک کند. در [۱۹] روشی برای استفاده از یک مدل توجه چند سطحی<sup>۲۵</sup> پیشنهاد می‌شود. در ابتدا سه اندازه میدان پذیرش معمولی برای نقشه‌های توجه چند سطحی انتخاب می‌شوند. سپس یادگیری چندسطحی برای استخراج ویژگی‌های متمایز از این مناطق محلی معرفی می‌گردند. این روش نگرش جدیدی در مورد چگونگی استفاده از فعال‌سازهای شبکه عصبی، برای تولید مناطق چند مقیاسی - که برای طبقه‌بندی ریزدانه‌ای مفید هستند - ارائه می‌دهد و شامل دو مرحله است: ۱- انتخاب نورهایی که حداکثر فعال‌سازی را در سه نقشه ویژگی انتخاب شده دارند. این نقشه‌ها خروجی مدل‌های شبکه عصبی کانولوشنی هستند که از قبل روی تصاویر اندازه کامل، آموزش داده شده‌اند. ۲- آموزش شبکه‌های ظریف با این مناطق چند مقیاسی ایجاد شده. هر منطقه متمایز شده را می‌توان به عنوان یکی از ویژگی‌ها در نظر گرفت. سپس این نتایج برای پیش‌بینی نهایی ادغام می‌شوند. در [۲۰] به یکی از مشکلات بازیابی تصویر ریزدانه‌ای

## ۲- بازیابی تصویر ریزدانه‌ای با استفاده از شبکه خودتوجهی مکانی و مکانیسم برجسته‌سازی

شکل ۱. شمای کلی شبکه خودتوجهی مکانی با مکانیسم برجسته‌سازی



در [۲۱] شبکه خودتوجهی مکانی<sup>۲۸</sup> پیشنهاد شده است که شامل دو جزء اصلی می‌باشد. ابتدا یک شبکه عصبی کانولوشنی به عنوان استخراج‌کننده ویژگی<sup>۲۹</sup> پیاده‌سازی می‌شود که ویژگی‌های اولیه را از تصاویر ورودی از طریق چندین لایه کانولوشن استخراج می‌کند. سپس ماژول خودتوجهی مکانی با استفاده از مکانیسم توجه، ویژگی‌های جدید را ذخیره می‌کند.

### ۲-۱- استخراج‌کننده ویژگی

اخیراً، برای وظایف پردازش تصویر، یک رویکرد مرسوم برای استخراج ویژگی‌های اولیه، استفاده از یک شبکه عصبی کانولوشنی از قبل آموزش دیده، به منظور بهره‌مندی از مقدار اولیه وزن معنادار است. چنین شبکه‌های عصبی کانولوشنی می‌توانند ویژگی‌های سطح بالا را از تصاویر استخراج کنند. برای مقایسه منصفانه با سایر روش‌های پیشرفته، از VGG-۱۶ از پیش آموزش دیده بر روی مجموعه داده ImageNet استفاده می‌شود.

برای استخراج اولیه، سه لایه آخر که کاملاً متصل هستند حذف می‌شوند. نقشه‌های ویژگی هر تصویر از مجموعه داده به عنوان  $X$  از خروجی لایه کانولوشنی نهایی گرفته می‌شود. این فرآیند به صورت زیر نشان داده شده است

$$F = VGG(X) \quad (1)$$

به طور خاص، استخراج‌کننده ویژگی، یک تصویر ورودی  $X$  را به یک نقشه ویژگی با ابعاد  $F \in R^{H \times W \times K}$  نگاشت می‌کند، که در آن  $H$ ,

در فرمول (۴)  $X$  نقشه ویژگی خودتوجهی مکانی تصویر ورودی و  $Y$  نقشه ویژگی خودتوجهی مکانی هر تصویر از پایگاه داده است. سپس فاصله‌های به‌دست‌آمده، که هرکدام نگاشتی به تصویری از پایگاه داده دارند، به صورت نزولی مرتب شده و پنج نتیجه برتر بازیابی می‌شود. خروجی سیستم بر اساس کلاسی که بیشترین احتمال را در بین این پنج نتیجه دارد، تعیین می‌گردد.

## ۲-۴- مکانیسم برجسته‌سازی

یکی از مشکلات روش استفاده از شبکه خودتوجهی مکانی آن است که تصویر ورودی، با ویژگی‌های با اهمیت کمتر بررسی می‌شود و ممکن است بخش‌های حاشیه‌ای در نتیجه نهایی عملکرد شبکه، تأثیرگذار باشند. در این مقاله، روش XRAI برجسته‌سازی<sup>۳۰</sup> پیشنهاد شده است. از آنجا که امکان‌پذیر نیست نتایج انتساب، به مجموعه خاصی از پارامترهای فوق یا کیفیت روش تقسیم‌بندی، بستگی داشته باشد، تصویر چندین بار با استفاده از مجموعه پارامترهای مختلف قطعه‌بندی می‌شود. به طور خاص، از یک پارامتر مقیاس در مجموعه [۵۰، ۱۰۰، ۲۵۰، ۵۰۰، ۱۲۰۰] استفاده شده و بخش‌های کوچکتر از بیست پیکسل نادیده گرفته می‌شود. یک پارامتر واحد، اجتماع بخش‌ها در کل تصویر را محاسبه می‌کند. بنابراین، اجتماع همه بخش‌ها مساحتی برابر با شش برابر مساحت تصویر ایجاد می‌کند و در نتیجه بخش‌های جداگانه به طور قابل توجهی همپوشانی دارند. مرزهای بخش معمولاً با لبه‌های تصویر همسو می‌شوند. برای استخراج نقشه‌های برجسته، مطلوب است که بخش‌ها شامل لبه‌ها باشند، زیرا تصاویر در دو طرف یک لبه نازک اغلب به یکدیگر مرتبط هستند. برای این منظور، ماسک‌های بخش‌ها پنج پیکسل گسترش می‌یابند تا مجموعه نهایی قطعات به دست آید.

### XRAI Algorithm:

Given Image  $I$ , model  $f$  and attribution method  $g$   
 Over-segment  $I$  to segments  $s \in S$   
 Get attribution map  $A = g(f, I)$   
 Let saliency mask  $M = 0$ , trajectory  $T = []$   
 while  $S \neq \emptyset$  and  $\text{area}(M) < \text{area}(I)$  do  
   for  $s \in S$   
     Compute gain:  $g_s = \sum_{i \in S \setminus M} \frac{A_i}{\text{area}(S \setminus M)}$   
   end for  
    $\hat{s} = \text{argmax}_s g_s$   
    $S = S \setminus \hat{s}$   
    $M = M \cup \hat{s}$   
   Add  $M$  to list  $T$   
end while  
return  $T$

## ۳- ارزیابی دقت و نتایج آزمایشگاهی

در تکمیل بهبود استفاده از شبکه خودتوجهی مکانی، از روش برجسته‌سازی XRAI استفاده شده است. مجموعه داده Stanford Dogs [۲۲] شامل ۲۰۵۸۰ تصویر از ۱۲۰ نژاد سگ از سراسر جهان است. این مجموعه داده با استفاده از تصاویر و حاشیه نویسی از ImageNet، برای طبقه بندی تصاویر ریز دانه ساخته شده است. این مجموعه داده یک مشکل چالش برانگیز داشت، زیرا برخی از نژادهای سگ ویژگی‌های تقریباً یکسانی دارند یا از نظر رنگ و سن متفاوت

$W$  و  $K$  نشان دهنده ارتفاع مکانی، عرض مکانی، تعداد کانال‌ها/کرانل حاوی کانال هستند.

## ۲-۲- مازول خودتوجهی مکانی

ماژول خودتوجهی مکانی از مکانیسم خودتوجهی پیشنهاد شده استفاده می‌کند که توجه محلی را از طریق یک تابع فعال‌ساز softmax جمع می‌کند. این ایده گسترش می‌یابد تا به موقعیت‌های پیکسل مکانی ویژگی‌های اصلی توجه و از تجمیع ویژگی‌ها برای به‌دست‌آوردن نقشه‌های ویژگی خودتوجهی مکانی استفاده شود. با توجه به نقشه‌های ویژگی اولیه  $F \in R^{H \times W \times K}$  به دست آمده از استخراج‌کننده ویژگی، ابتدا سه نقشه ویژگی جدید  $A, B, C$  با استفاده از کانولوشن  $1 \times 1$  تولید می‌شود.  $\{A, B, C\} \in R^{H \times W \times K}$  همان ابعاد فضای  $F$  را داراست. سپس  $A$  و  $B$  و  $C$  به  $R^{N \times K}$  تغییر شکل می‌یابد، که در آن  $N = H \times W$  تعداد پیکسل‌ها است. سپس، ضرب عناصر بین  $A$  و ترانپاده  $B$  محاسبه می‌شود. تابع فعال‌ساز softmax از نظر مکانی برای محاسبه نقشه خودتوجهی مکانی اعمال می‌شود  $S \in R^{N \times K}$  که:

$$S_{ij} = \frac{\exp(A_i \otimes B_j)}{\sum_{i=1}^N \exp(A_i \otimes B_j)} \quad (2)$$

در فرمول (۲)  $\otimes$  نشان‌دهنده ضرب عنصر به عنصر است.  $S_{ij}$  نشان می‌دهد که چگونه شبکه تأثیر  $i$ مین موقعیت مکانی را بر موقعیت مکانی  $j$ مین اندازه گیری می‌کند. از این رو، بازنمایی ویژگی‌های مرتبط‌تر بین  $A$  و  $B$  منجر به همبستگی معنی‌دار و غنی‌تر بین آنها می‌شود و بالعکس. برای تقویت موقعیت‌های مکانی، ضرب عناصر بین  $S \in R^{N \times K}$  و  $C \in R^{N \times K}$  انجام می‌شود و نتایج به  $R^{H \times W \times K}$  تغییر شکل داده می‌شود.

در نهایت، یک مکانیسم تجمیع ویژگی برای بررسی تأثیر مناطق خودتوجهی مکانی در همه موقعیت‌ها، در نقشه ویژگی اصلی از طریق فرمول (۳) پیاده‌سازی می‌شود:

$$H_j = \sum_{i=1}^N (S_{ij} C_i) \oplus F_j \quad (3)$$

می‌توان از فرمول (۳) استنباط کرد که ویژگی‌های به‌دست‌آمده توسط  $H_j$  نشان‌دهنده یک تجمع کلی از نمای زمینه‌ای بر اساس نقشه‌های خودتوجهی مکانی است. مجموعه این ویژگی‌ها به عنوان یک پایگاه داده ذخیره می‌شوند.

## ۲-۳- بازیابی تصویر

در این بخش یک تصویر به عنوان ورودی به شبکه داده می‌شود و طبق معادلات (۱)، (۲) و (۳) نقشه ویژگی‌های خودتوجهی مکانی آن به دست می‌آید. سپس با نقشه‌های ویژگی ذخیره شده در قسمت ۲-۲ و با استفاده از معادله زیر مقایسه می‌شوند:

$$D(X, Y) = \sqrt{\sum_{i=1}^{N=20580} (X_i - Y_i)^2} \quad (4)$$



EfficientNet-B۰ [۲۸]	۶۱٪
PC [۲۹]	۶۱٪
SSA [۳۰]	۸۶٪
روش پیشنهادی	۸۸٪

جدول ۲ بیان‌کننده مقایسه نتایج عملکرد روش‌های مختلف در بازیابی تصویر ریزدانه‌ای بر روی مجموعه داده Stanford Dogs می‌باشد. در روش EfficientNet-B۰ از کانولوشن‌های قابل تفکیک عمیق استفاده شده و در میان سایر روش‌ها از عملکرد ضعیف‌تری برخوردار است. در روش PC از آموزش با تقسیم‌بندی‌های پیچیده، در جهت ریزدانه‌ای کردن بازیابی تصویر استفاده شده‌است که نسبت به روش پیشین ۰/۷ درصد بهبود داشته‌است. روش PDFR، از فیلترهای متمایز و آشکارسازی استفاده می‌کند که نسبت به روش قبلی ۱۰ درصد بهبود ایجاد کرده‌است. در روش‌های FCAN، PC-DenseNet، HDWE و ۱۱۶ نیز از شبکه‌های کانولوشنی متعددی استفاده شده که نتیجه آن‌ها نسبت به موارد قبلی بین ۸ الی ۱۳ درصد افزایش امتیاز صورت گرفته‌است. روش SSA که در این مقاله نیز بررسی شد نسبت به سایر روش‌ها کارآمدتر است ولی روش پیشنهادی عملکردی بهتری نسبت به تمامی روش‌های ذکر شده دارد.

#### ۴- نتیجه‌گیری

در این مقاله یک روش بازیابی تصویر ریزدانه‌ای با کمک شبکه خودتوجهی مکانی و مکانیسم توجه پیشنهاد شده‌است. شبکه خودتوجهی مکانی با استخراج ویژگی و استفاده از مکانیسم توجه، ویژگی‌های جدید تصاویر را ذخیره می‌کند و مکانیسم برجسته‌سازی با حذف حاشیه‌های کم‌اهمیت تصویر ورودی، عملکرد سیستم را ارتقا می‌بخشد. نتایج نشان می‌دهند که روش پیشنهادی در بازیابی تصویر ریزدانه‌ای عملکرد مطلوبی دارد.

#### مراجع

- [۱] Li, J., Ling, Z., Niu, L., & Zhang, L. (۲۰۲۲). Zero-shot sketch-based image retrieval with structure-aware asymmetric disentanglement. In Computer Vision and Image Understanding (Vol. ۲۱۸, p. ۱۰۳۴۱۲). Elsevier BV. <https://doi.org/10.1016/j.cviu.2022.103412>
- [۲] Wang, Y., Huang, F., Zhang, Y., Feng, R., Zhang, T., & Fan, W. (۲۰۲۰). Deep cascaded cross-modal correlation learning for fine-grained sketch-based image retrieval. In Pattern Recognition (Vol. ۱۰۰, p. ۱۰۷۱۴۸). Elsevier BV. <https://doi.org/10.1016/j.patcog.2019.107148>
- [۳] Bressan, R. S., Bugatti, P. H., & Saito, P. T. M. (۲۰۲۲). Optimum-path forest and active learning approaches for content-based medical image retrieval. In Optimum-Path Forest (pp. ۹۵-۱۰۷). Elsevier. <https://doi.org/10.1016/b978-0-12-822688-9.00128>
- [۴] Rahman, A., Winarko, E., & Mustofa, K. (۲۰۲۲). Product image retrieval using category-aware siamese convolutional neural network feature. In Journal of King Saud University - Computer and Information Sciences (Vol. ۳۴, Issue ۶, pp. ۲۶۸۰-۲۶۸۷). Elsevier BV. <https://doi.org/10.1016/j.jksuci.2022.03.005>
- [۵] Jiang, Y., Du, J., Xue, Z., & Li, A. (۲۰۲۲). Cross-Media Retrieval of Scientific and Technological Information Based on Multi-Feature Fusion. In Neurocomputing. Elsevier BV. <https://doi.org/10.1016/j.neucom.2022.06.061>
- [۶] Niu, D., Zhao, X., Lin, X., & Zhang, C. (۲۰۲۰). A novel image retrieval method based on multi-features fusion. In Signal Processing: Image Communication (Vol. ۸۷, p. ۱۱۵۹۱۱). Elsevier BV. <https://doi.org/10.1016/j.image.2020.115911>

هستند. در آزمایش پیش رو، برای شروع، تصویر ورودی به شبکه داده می‌شود و طبق معادلات (۱)، (۲) و (۳) نقشه ویژگی‌های خودتوجهی مکانی آن به دست می‌آید. سپس با نقشه‌های ویژگی ذخیره شده در قسمت ۲ و با استفاده از معادله (۴) مقایسه می‌شوند. سپس فاصله‌های به‌دست‌آمده، که هرکدام نگاشتی به تصویری از پایگاه داده دارند، به صورت نزولی مرتب شده و ۵ نتیجه برتر بازیابی می‌شود. امتیاز سیستم بر اساس کلاسی که بیشترین احتمال را در بین این ۵ نتیجه دارد، تعیین می‌گردد.

$$Score = \frac{True Prediction}{All Prediction} \quad (6)$$

در این آزمایش از یک کامپیوتر شخصی مدل Macbook Air با پردازنده M1 و حافظه اصلی ۸ گیگابایت استفاده شده است. محیط توسعه و زبان مورد استفاده python می‌باشد.



شکل ۲۱. نمونه‌ای از تصاویر مجموعه داده [۲۳]

جدول ۱. نتایج آزمایش روی داده‌های آزمایشی

Input	Outputs						Score
							۸۰٪
Lhasa	Lhasa	Lhasa	Lhasa	Lhasa	Lhasa	silky_terrier	
							۸۰٪
Irish setter	Irish setter	Irish setter	Irish setter	Irish setter	Irish setter	Irish setter	
							۱۰۰٪
Pomeranian	Pomeranian	Pomeranian	Pomeranian	Pomeranian	Pomeranian	Pomeranian	
							۶۰٪
Afghan_hound	bloodhound	Afghan_hound	bloodhound	Afghan_hound	Afghan_hound	Afghan_hound	
							۱۰۰٪
otterhound	otterhound	otterhound	otterhound	otterhound	otterhound	otterhound	
							۱۰۰٪
Afghan_hound	Afghan_hound	Afghan_hound	Afghan_hound	Afghan_hound	Afghan_hound	Afghan_hound	
							۶۰٪
Pekinese	Pekinese	Shih-Tzu	Shih-Tzu	Pekinese	Pekinese	Pekinese	
							۱۰۰٪
Bernese_mountain	Bernese_mountain	Bernese_mountain	Bernese_mountain	Bernese_mountain	Bernese_mountain	Bernese_mountain	
							۱۰۰٪
chow	chow	chow	chow	chow	chow	chow	
							۱۰۰٪
clumber	clumber	clumber	clumber	clumber	clumber	clumber	

جدول ۲. نتایج آزمایشگاهی

Method	Score
FCAN [۲۴]	۸۴/۵٪
PDFR [۲۵]	۷۷/۹٪
PC-DenseNet-۱۶۱ [۲۶]	۸۳/۶٪
HDWE [۲۷]	۷۹/۶٪

- [۲۹] Ji, J., Guo, Y., Yang, Z., Zhang, T., & Lu, X. (۲۰۲۱). Multi-level dictionary learning for fine-grained images categorization with attention model. *Neurocomputing*, ۴۵۳, ۴۰۳-۴۱۲. <https://doi.org/10.1016/j.neucom.2020.07.147>
- [۳۰] Zeng, X., Zhang, Y., Wang, X., Chen, K., Li, D., & Yang, W. (۲۰۱۹). Fine Grained Image Retrieval via Piecewise Cross Entropy loss. *Image and Vision Computing*. <https://doi.org/10.1016/j.imavis.2019.10.006>
- [۳۱] Baffour, A. A., Qin, Z., Wang, Y., Qin, Z., & Choo, K. -K. R. (۲۰۲۱). Spatial self-attention network with self-attention distillation for fine-grained image recognition. *Journal of Visual Communication and Image Representation*, ۸۱, ۱۰۳۳۶۸. <https://doi.org/10.1016/j.jvcir.2021.103368>
- [۳۲] Khosla, A., Jayadevaprakash, N., Yao, B., & Fei-Fei, L. (۲۰۱۲). Novel Dataset for Fine-Grained Image Categorization : Stanford Dogs.
- [۳۳] Khosla, A., Jayadevaprakash, N., Yao, B., & Fei-Fei, L. (۲۰۱۲). Novel Dataset for Fine-Grained Image Categorization : Stanford Dogs.
- [۳۴] Liu, X., Xia, T., Wang, J., & Lin, Y. (۲۰۱۶). Fully convolutional attention localization networks: Efficient attention localization for fine-grained recognition. *arXiv preprint arXiv:1607.06765*, ۱(۲), ۴.
- [۳۵] X. Zhang, H. Xiong, W. Zhou, W. Lin and Q. Tian, "Picking Deep Filter Responses for Fine-Grained Image Recognition," ۲۰۱۶ IEEE Conference on Computer Vision and Pattern Recognition (CVPR), ۲۰۱۶, pp. ۱۱۳۴-۱۱۴۲, doi: 10.1109/CVPR.2016.128.
- [۳۶] Dubey, A., Gupta, O., Guo, P., Raskar, R., Farrell, R., Naik, N. (۲۰۱۸). Pairwise Confusion for Fine-Grained Visual Classification. In: Ferrari, V., Hebert, M., Sminchisescu, C., Weiss, Y. (eds) *Computer Vision – ECCV ۲۰۱۸*, ECCV ۲۰۱۸, Lecture Notes in Computer Science(), vol ۱۱۲۱۶. Springer, Cham. [https://doi.org/10.1007/978-3-030-12588-8\\_5](https://doi.org/10.1007/978-3-030-12588-8_5)
- [۳۷] Yu, J., Huang, Y., Gbur, G., Wang, F., & Cai, Y. (۲۰۱۹). Enhanced backscatter of vortex beams in double-pass optical links with atmospheric turbulence. In *Journal of Quantitative Spectroscopy and Radiative Transfer* (Vol. ۲۲۸, pp. ۱-۱۰). Elsevier BV. <https://doi.org/10.1016/j.jqsrt.2019.02.021>
- [۳۸] D. Haase and M. Amthor, "Rethinking Depthwise Separable Convolutions: How Intra-Kernel Correlations Lead to Improved MobileNets," ۲۰۲۰ IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), ۲۰۲۰, pp. ۱۴۵۸۸-۱۴۵۹۷, doi: 10.1109/CVPR42600.2020.01461.
- [۳۹] Dubey, A., Gupta, O., Guo, P., Raskar, R., Farrell, R., & Naik, N. (۲۰۱۷). Training with confusion for fine-grained visual classification. *CoRR*
- [۴۰] Baffour, A. A., Qin, Z., Wang, Y., Qin, Z., & Choo, K. -K. R. (۲۰۲۱). Spatial self-attention network with self-attention distillation for fine-grained image recognition. *Journal of Visual Communication and Image Representation*, ۸۱, ۱۰۳۳۶۸. <https://doi.org/10.1016/j.jvcir.2021.103368> زیرنویس‌ها
- [۷] Zhang, L., & Wu, X. (۲۰۲۲). Multi-task framework based on feature separation and reconstruction for cross-modal retrieval. In *Pattern Recognition* (Vol. ۱۲۲, p. ۱۰۸۲۱۷). Elsevier BV. <https://doi.org/10.1016/j.patcog.2021.108217>
- [۸] Zhang, K., Qi, S., Cai, J., Zhao, D., Yu, T., Yue, Y., Yao, Y., & Qian, W. (۲۰۲۲). Content-based image retrieval with a Convolutional Siamese Neural Network: Distinguishing lung cancer and tuberculosis in CT images. In *Computers in Biology and Medicine* (Vol. ۱۴۰, p. ۱۰۵۰۹۶). Elsevier BV. <https://doi.org/10.1016/j.combiomed.2021.105096>
- [۹] Yelchuri, R., Dash, J. K., Singh, P., Mahapatro, A., & Panigrahi, S. (۲۰۲۲). Exploiting deep and hand-crafted features for texture image retrieval using class membership. In *Pattern Recognition Letters* (Vol. ۱۶۰, pp. ۱۶۳-۱۷۱). Elsevier BV. <https://doi.org/10.1016/j.patrec.2022.06.017>
- [۱۰] Chen, Y., Zhang, Z., Wang, Y., Zhang, Y., Feng, R., Zhang, T., & Fan, W. (۲۰۲۲). AE-Net: Fine-grained sketch-based image retrieval via attention-enhanced network. *Pattern Recognition*, ۱۲۲, ۱۰۸۲۹۱. <https://doi.org/10.1016/j.patcog.2021.108291>
- [۱۱] Zeng, X., Liu, S., Wang, X., Zhang, Y., Chen, K., & Li, D. (۲۰۲۱). Hard Decorrelated Centralized Loss for fine-grained image retrieval. *Neurocomputing*, ۴۵۳, ۲۶-۳۷. <https://doi.org/10.1016/j.neucom.2021.04.030>
- [۱۲] Shan, W., Huang, D., Wang, J., Zou, F., & Li, S. (۲۰۲۲). Self-Attention based fine-grained cross-media hybrid network. In *Pattern Recognition* (Vol. ۱۳۰, p. ۱۰۸۷۴۸). Elsevier BV. <https://doi.org/10.1016/j.patcog.2022.108748>
- [۱۳] Guo, C., Lin, Y., Chen, S., Zeng, Z., Shao, M., & Li, S. (۲۰۲۲). From the whole to detail: Progressively sampling discriminative parts for fine-grained recognition. *Knowledge-Based Systems*, ۲۳۵, ۱۰۷۶۵۱. <https://doi.org/10.1016/j.knosys.2021.107651>
- [۱۴] Wei, X.-S., Xie, C.-W., Wu, J., & Shen, C. (۲۰۱۸). Mask-CNN: Localizing parts and selecting descriptors for fine-grained bird species categorization. In *Pattern Recognition* (Vol. ۷۶, pp. ۷۰۴-۷۱۴). Elsevier BV. <https://doi.org/10.1016/j.patcog.2017.10.002>
- [۱۵] Jin, X., Deng, Q., Lou, H., Li, X., & Xiao, C. (۲۰۲۲). Fine-grained Regression for Image Aesthetic Scoring. In *Cognitive Robotics*. Elsevier BV. <https://doi.org/10.1016/j.cogr.2022.07.002>
- [۱۶] Ha, Y., Du, Z., & Tian, J. (۲۰۲۲). Fine-grained interactive attention learning for semi-supervised white blood cell classification. *Biomedical Signal Processing and Control*, ۷۵, ۱۰۳۶۱۱. <https://doi.org/10.1016/j.bspc.2022.103611>
- [۱۷] Zhao, Q., Wang, X., Lyu, S., Liu, B., & Yang, Y. (۲۰۲۲). A feature consistency driven attention erasing network for fine-grained image retrieval. *Pattern Recognition*, ۱۲۸, ۱۰۸۶۱۸. <https://doi.org/10.1016/j.patcog.2022.108618>
- [۱۸] Yan, T., Shi, J., Li, H., Luo, Z., & Wang, Z. (۲۰۲۲). Discriminative information restoration and extraction for weakly supervised low-resolution fine-grained image recognition. *Pattern Recognition*, ۱۲۷, ۱۰۸۶۲۹. <https://doi.org/10.1016/j.patcog.2022.108629>

۱۶ Image Aesthetic Scoring

۱۷ Mean Square Errors

۱۸ Segmented Mean Square Errors

۱۹ Hash Learning Method

۲۰ Feature Consistency Driven Attention Erasing Network:

FCAENet

۲۱ Selective Region Erasing Module: SREM

۲۲ Enhancing Space Relation Loss: ESRL

۲۳ Fine-Grained discriminative Information Restoration: FDR

۲۴ Semantic Relation Distillation Loss: SRD-Loss

۲۵ Multi-level Attention Model

۲۶ Cross Entropy Loss

۲۷ Piecewise Cross Entropy loss

۲۸ Spatial Self-Attention Network (SSA.Net)

۲۹ Feature Extractor: FE

۳۰ Saliency

۱ Text Based Image Retrieval

۲ Content Based Image Retrieval

۳ Fine-Grained Content Based Image Retrieval

۴ Zero-Shot Sketch based Image Retrieval (ZS-SBIR)

۵ Asymmetric Disentanglement

۶ Deep Cascaded Cross-modal Ranking Model

۷ Siamese Convolutional Network

۸ Multi-feature Fusion based Cross-Media Retrieval

۹ Adversarial Learning

۱۰ Uniform local binary patterns

۱۱ multi-task framework based on feature separation and reconstruction

۱۲ Transformator

۱۳ Self-Attention Network

۱۴ Cross-Media

۱۵ Mask-Convolutional Neural Network