



دانشگاه صنعتی امیر کبیر  
(پلی تکنیک تهران)

درس سمینار رشته هوش مصنوعی

# تولید خودکار شرح بر تصاویر با استفاده از شبکه‌های عصبی کانولوشنی عمیق و بازگشتی

Automatic Image Captioning Using Deep Convolutional and Recurrent Neural Networks

استاد راهنما

دکتر صفابخش

پژوهش گر

احمد اسدی

۹۴۱۳۱۰۹۱

اردیبهشت ماه ۱۳۹۵

## چکیده

با توجه به افزایش چشمگیر تعداد تصاویر مورد استفاده کاربران در فضاهای مجازی و همین‌طور با در نظر گرفتن گرایش روزافزون کاربران به ذخیره‌سازی تصاویر در رایانه‌های شخصی، مساله مدیریت این تصاویر و یافتن تصاویر خاص بین مجموعه تصاویر موجود، به یکی از مسائل مهم و پرکاربرد در زمینه بینایی ماشین تبدیل شده است. گام اساسی در این راستا، دستیابی به سامانه‌ای است که قادر به تولید خودکار شرح برای تصاویر باشد. شرح این تصاویر که در قالب جملات زبان طبیعی ارائه می‌شود باید علاوه بر سازگاری با موضوع تصویر و توصیف صحیح صحنه، به لحاظ دستور زبان و معنا صحیح و کامل باشد.

# فهرست مطالب

۱	فصل اول مقدمات
۲	۱-۱ مقدمه
۳	۲-۱ تعریف مساله
۴	۳-۱ درک صحنه
۵	۱-۳-۱ پژوهش‌های انجام شده در زمینه درک صحنه توسط مغز انسان
۷	۲-۳-۱ نتایج به دست آمده از آزمایشات
۹	۴-۱ جمع‌بندی
۱۱	فصل دوم درک صحنه
۱۲	۱-۲ درک صحنه
۱۲	۲-۲ روش‌های مختلف موجود
۱۳	۳-۲ روش‌های مبتنی بر مدل‌های گرافی احتمالی
۱۳	۱-۳-۲ استفاده از مدل میدان تصادفی مارکف
۱۶	۲-۳-۲ استفاده از مدل میدان تصادفی شرطی
۱۸	۳-۳-۲ استفاده از سایر مدل‌های گرافی احتمالی
۲۵	۴-۲ روش‌های مبتنی بر شبکه‌های عصبی کانولوشنی عمیق
۲۵	۱-۴-۲ اختصاص معنا به قطعه‌های مختلف تصویر
۲۶	۲-۴-۲ ناحیه‌بندی عمیق تصاویر به منظور نگاشت دو طرفه جملات و تصاویر
۳۱	۵-۲ جمع‌بندی
۳۵	فصل سوم تولید جمله
۳۶	۱-۳ تولید جمله
۳۶	۲-۳ کاربردها
۳۷	۳-۳ روش تولید زبان طبیعی
۳۸	۴-۳ روش نزدیک‌ترین همسایه
۴۲	۵-۳ استفاده از قالب‌های آماده زبانی
۴۴	۶-۳ روش‌های مبتنی بر شبکه‌های عصبی بازگشتی
۴۶	۱-۶-۳ شبکه عصبی بازگشتی ضربی
۴۸	۲-۶-۳ تولید جمله با مفهوم مشخص
۵۳	۷-۳ جمع‌بندی
۵۶	۴ فصل چهارم تولید شرح متناظر صحنه با استفاده از یادگیری عمیق
۵۷	۱-۴ تولید شرح بر تصاویر با استفاده از یادگیری عمیق

۵۸	۵	فصل پنجم تولید شرح متناظر صحنه با استفاده از روش‌های مبتنی بر نقطه توجه
۵۹	۱-۵	تولید شرح بر تصاویر با استفاده از روش‌های مبتنی بر نقطه توجه
۵۹	۲-۵	روش‌های مبتنی بر نقطه توجه در حوزه ترجمه ماشینی
۶۰	۱-۲-۵	انکودر
۶۱	۲-۲-۵	دیکودر
۶۱	۳-۲-۵	ایده اصلی استفاده از نقطه توجه
۶۲	۴-۲-۵	دیکودر در روش مبتنی بر نقطه توجه
۶۳	۵-۲-۵	انکودر در روش مبتنی بر نقطه توجه
۶۴	۶	فصل ششم تولید شرح متناظر صحنه با استفاده از روش‌های مبتنی بر حافظه فعال
۶۵	۱-۶	تولید جمله
۶۶	۷	فصل هفتم جمع‌بندی و نتیجه‌گیری
۶۷	۱-۷	جمع‌بندی و نتیجه‌گیری

## فهرست تصاویر

۶	نمونه توصیف‌های افراد برای تصاویر . . . . .	۱
۶	ساختار مطلوب اطلاعات استخراج شده از تصاویر [۱] . . . . .	۲
۷	تصاویر دنیای واقعی مورد استفاده در آزمایشات [۱] . . . . .	۳
۸	نمودار مقایسه‌ای عملکرد مغز در درک صحنه . . . . .	۴
۹	نمونه‌ای از نتایج بهدست‌آمده از آزمایشات [۱] . . . . .	۵
۱۴	نگاشت تصویر به فضای معنایی . . . . .	۶
۱۵	مدل میدان تصادفی مارکف در درک صحنه . . . . .	۷
۱۷	مدل سلسله‌مراتبی میدان تصادفی شرطی در درک صحنه . . . . .	۸
۱۸	مدل گرافی احتمالی مورد استفاده در پژوهش [۲] . . . . .	۹
۲۲	نمونه تصاویر موجود در مجموعه‌داده مورد استفاده [۲] . . . . .	۱۰
۲۳	ماتریس درهم‌ریختگی مدل کامل ارائه شده در [۲] . . . . .	۱۱
۲۳	نتیجه مقایسه مدل‌های مختلف در [۲] . . . . .	۱۲
۲۴	نتایج نهایی بهدست آمده از مدل بر روی تصاویر. [۲] . . . . .	۱۳
۲۷	طرح‌واره عملکرد روش RCNN . . . . .	۱۴
۲۹	نتایج عملکرد اهداف تعریف شده در روش RCNN برای همترازسازی تصاویر و جملات . . . . .	۱۵
۳۰	نتایج نهایی روش RCNN . . . . .	۱۶
۳۱	نتایج حاصل از جستجوی جملات در روش RCNN . . . . .	۱۷
۳۹	نتایج کیفی اختصاص جملات و تصاویر به یکدیگر با استفاده از روش نزدیک‌ترین همسایه [۳] . . . . .	۱۸
۴۱	نتایج برنامه‌سازی خطی صحیح در مقایسه با جملات تولید شده انسان [۴] . . . . .	۱۹
۴۱	مواردی از خروجی برنامه‌سازی خطی صحیح که نسبت به جملات انسان، برتری دارد [۴] . . . . .	۲۰
۴۲	نتایج برنامه‌سازی خطی صحیح که به لحظه‌های مختلف دچار مشکل شده‌اند [۴] . . . . .	۲۱
۴۴	نمونه‌های صحیح از جملات تولید شده توسط قالب‌های آماده زبانی [۵] . . . . .	۲۲
۴۴	نمونه‌های اشتباه تولید شده توسط قالب‌های آماده زبانی [۵] . . . . .	۲۳
۴۷	طرح‌واره شبکه عصبی بازگشتی ضربی [۶] . . . . .	۲۴
۴۸	هم‌ترازسازی تصویر و جمله [۷] . . . . .	۲۵
۴۹	ارتباط بین نواحی مختلف یک تصویر و عبارات جمله [۷] . . . . .	۲۶
۵۰	طرح‌واره شبکه عصبی بازگشتی دو طرفه [۷] . . . . .	۲۷
۵۱	انتساب نواحی مختلف تصویر به عبارات زبانی [۷] . . . . .	۲۸
۵۲	طرح‌واره شبکه عصبی بازگشتی ارائه شده برا تولید جمله [۷] . . . . .	۲۹
۵۳	نتایج رتبه‌بندی عبارات زبانی برای نواحی تصویر [۷] . . . . .	۳۰
۵۳	نتایج تولید جمله برای تصاویر در [۷] . . . . .	۳۱
۶۰	ساختار کلی چارچوب کاری انکوادر-دیکوادر . . . . .	۳۲

٦١	ساختار دیکودر مورد استفاده در چارچوب کاری [۲۸]	۳۳
٦٣	ساختار کلی یک شبکه عصبی بازگشته دوطرفه . . . . .	۳۴

## فهرست جداول

۳۹ . . . . .	نتایج استفاده از روش نزدیک‌ترین همسایه در اختصاص جملات و تصاویر [۳]	۱
۴۴ . . . . .	نتایج معیارهای ROUGE و BLUE در حالات مختلف [۵]	۲
۵۲ . . . . .	نتایج معیار BLUE برای روش ارائه شده در [۷] در مقایسه با دو روش دیگر	۳

# ١ فصل اول

## مقدمات

به دنبال پیشرفت تکنولوژی در ساخت دوربین‌های عکاسی و ورود دوربین‌های نیمه‌خودکار و خودکار به بازار، تعداد زیادی از کاربران سیستم‌های رایانه‌ای به استفاده از این تکنولوژی در ثبت تصاویر مورد علاقه خود جذب شده‌اند. دقیق و کیفیت مطلوب تصویربرداری از یک سو و سهولت استفاده از دوربین از سوی دیگر، باعث شده‌اند تعداد تصاویر ثبت شده توسط کاربران به طور روزافزون افزایش یابد؛ به‌طوری‌که امروزه اغلب کاربران، تعداد بی‌شماری از این تصاویر را در گوشی‌های تلفن همراه، تبلت‌ها و رایانه‌های شخصی خود نگه‌داری می‌کنند. از جمله مشکلاتی که در اثر ایجاد این حجم وسیع از تصاویر بوجود آمده، مشکل مدیریت این تصاویر و یافتن تصاویر خاص بین مجموعه بزرگی از تصاویر موجود، است.

برای دست‌یابی به سامانه‌ای که بتواند تعداد زیادی از تصاویر موجود را مدیریت نماید، ابتدا باید صحنه موجود در تصویر را به درستی درک کرد. درک صحیح از صحنه، عبارت است از بیان تصویر به نحوی که اطلاعات کلی موجود و هدف اصلی تصویر، واضح و مشخص باشد. این بیان می‌تواند شامل اجسام موجود در تصویر، رابطه مکانی بین اجسام، فعالیت به تصویر کشیده شده، شرایط محیطی موثر بر صحنه و مواردی از این دست باشد. از طرفی باید به نحوی محتوای تصاویر را بیان کرد که بتوان عملیات جستجو را بر اساس مدل بیان شده تصاویر انجام داد. در این صورت به‌ازای هر تصویر، یک نمونه از مدل مطابق با تصویر ایجاد و ذخیره خواهد شد. پرس‌وجوی<sup>۱</sup> کاربر به فضای مدل، نگاشت شده و تصویر معادل با مدل استخراج شده، به عنوان نتیجه جستجو نمایش داده می‌شود. علاوه بر این، مساله مدیریت تصاویر، به مساله مدیریت مدل‌های موجود کاهش داده می‌شود.

تولید شرح کلی بر تصاویر،<sup>۲</sup> بیان مناسبی از صحنه موجود در تصویر را ارائه می‌دهد. شرح تولید شده بر تصاویر، در قالب مجموعه‌ای از جملات زبان طبیعی<sup>۳</sup> ارائه می‌شود که عموماً بیان‌گر اجسام موجود در صحنه، ارتباطات مکانی بین اجسام و اطلاعات مشخص دیگر است که در هر پژوهش می‌تواند متفاوت باشد. بنابراین، دست‌یابی به سامانه‌ای که قادر به تولید خودکار شرح کلی بر تصاویر باشد، اساسی‌ترین گام در راستای تولید نرم‌افزارهای مدیریت تصاویر است.

یکی از اولین ایده‌های مطرح شده در این زمینه، با الهام از پژوهش‌های صورت گرفته در زمینه ترجمه ماشین<sup>۴</sup> به وجود آمده است که با هدف ترجمه جملات یک زبان به زبان دیگر به طور خودکار، انجام شده‌اند. در این راستا،

<sup>۱</sup>Query

<sup>۲</sup>Holistic Image Caption

<sup>۳</sup>Natural Language Sentences

<sup>۴</sup>Machine Translation

یک جمله از زبان مبدا<sup>۵</sup>، با روش‌های مختلف، تبدیل به یک بردار ویژگی<sup>۶</sup> می‌شود که مشخصه‌های اصلی جمله اولیه را نمایش می‌دهد. سپس بردار ویژگی حاصل، با اعمال روش‌های گوناگون دیگری، تبدیل به یک جمله از زبان مقصد<sup>۷</sup> میگردد که در آن تمام ویژگی‌های موجود در بردار ویژگی بیان شده‌اند. با توجه به فرایند مذکور، اگر به جای جمله زبان مبدا، یک تصویر به بردار ویژگی تبدیل شده و سپس با استفاده از روش‌های موجود قبلی، بردار ویژگی حاصل به جمله زبان مقصد ترجمه شود، جمله‌ای معادل با تصویر ورودی به دست خواهد آمد که بیان‌گر محتوای به تصویر کشیده شده در تصویر ورودی است [۸].

شرح خودکار تصاویر، توجه پژوهش‌گران بسیار زیادی را به خود جلب کرده است و فعالیت‌های متنوع و متعددی در این راستا انجام شده‌اند. علی‌رغم وجود پژوهش‌های فراوان و متفاوت، می‌توان یک بستر کلی برای تمام فعالیت‌های موجود در این زمینه ارائه داد. بر این مبنای، فرایند کلی که در عموم پژوهش‌های انجام‌شده، پی‌گرفته شده‌است، از دو بخش اساسی تشکیل می‌شود.

۱. بازنمایی تصاویر، با استفاده از بردار ویژگی

۲. تبدیل بردار ویژگی به دست آمده به جملات صحیح زبانی

## ۲-۱ تعریف مساله

در این پژوهه قصد داریم سامانه‌ای ارائه دهیم که قادر به تولید شرح کوتاه بر تصاویر باشد. دو دیدگاه اساسی در دست‌یابی به چنین سامانه‌ای مطرح است.

۱. یافتن نقاط توجه<sup>۸</sup> در تصاویر و تولید جملات توصیف‌کننده اجسام مستقر در این نقاط به طوری که توصیف جسم مستقر در نقطه توجه و اجسام مرتبط با آن در جملات تولیدی، وجود داشته باشد.

۲. تولید شرح جامع بر تصاویر به طوری که تمام اجسام موجود در صحنه به همراه روابط موجود بین آن‌ها توصیف شوند.

شرح کوتاه تولید شده در این پژوهه، به معنی تولید جملاتی است که مستقیماً به توصیف صحنه، اجسام موجود در صحنه و روابط بین آنها می‌پردازند. به طور کلی، دو چالش عمده در این پژوهش مورد توجه قرار خواهد گرفت:

۱. توصیف صحنه باید دقیق باشد؛ به این معنی که اجسام موجود در صحنه باید به طور دقیق از هم تفکیک شده و دسته‌بندی شوند. تصویر توصیف شده باید در قالب مناسبی بازنمایی شود که بتوان به راحتی از آن برای تولید جمله استفاده نمود.

۲. جملات تولید شده برای شرح تصویر باید به لحاظ دستور زبان، املاء و معنا صحیح بوده و با تصویر مرتبط خود سازگار باشند و آن را به درستی و دقیق شرح دهند.

<sup>۵</sup>Source Language

<sup>۶</sup>Feature Vector

<sup>۷</sup>Destination Language

<sup>۸</sup>Attention Points

## ۳-۱ درک صحنه

مساله درک صحنه، یکی از چالش‌های بزرگ و قدیمی مطرح در زمینه بینایی ماشین است. در گذشته، هدف اغلب پژوهش‌گران از طرح این مساله، توصیف صحنه موجود در تصویر با دیدن لحظه‌ای تصویر بوده است؛ اگرچه امروزه، تعریف این مساله دچار تغییر شده است.

به طور کل نمی‌توان تعریف جامع و شاملی برای درک صحنه ارائه داد. اگرچه تعاریفی عمومی ارائه شده‌اند که کلیات این مفهوم را توضیح می‌دهند. پژوهش‌گران در این زمینه هریک سعی در ارائه تعریفی برای این مفهوم دارند که برای کاربرد مورد نظر خود کافی و مفید باشد. به عنوان مثال، یکی از جدیدترین تعاریف برای درک صحنه در پژوهش [۹] ارائه شده است:

\* «توانایی تحلیل بصری یک صحنه برای پاسخ‌دادن به سوالاتی مانند سوالات زیر:

- چه اتفاقی در حال رخ دادن است؟

- چرا این اتفاق در حال رخ دادن است؟

- اتفاق بعدی که رخ خواهد داد، چیست؟»

به طور کل می‌توان درک صحنه را چنین معنا کرد:

\* درک صحنه، فرایندی است که طی آن، یک سامانه رایانه‌ای با استفاده از الگوریتم‌های موجود، اطلاعات بصری نهفته در تصاویر را استخراج کرده و در قالب مناسبی بازنمایی<sup>۹</sup> کند به طوری که این اطلاعات برای توصیف صحنه کافی و مفید باشد.

با این تعریف، اگرچه مفهوم درک صحنه کمی روشن می‌شود اما نکات مبهمی مانند این که چه نوع اطلاعاتی از تصویر استخراج شود، نیاز به توضیح و تفسیر بیشتری دارند. در تمام پژوهش‌های موجود در زمینه درک صحنه، که تعدادی از آن‌ها را در فصل بعدی مورد بررسی قرار خواهیم داد، تعریف اتخاذ شده برای درک صحنه، همین تعریف است با این تفاوت که اطلاعات مورد نیاز برای استخراج، در هر پژوهش، بسته به کاربرد تعریف می‌شود. موارد مختلفی که به عنوان اطلاعات لازم برای درک و توصیف صحنه، در پژوهش‌ها به چشم می‌خورد عموماً شامل موارد زیر هستند:

۱. دسته صحنه<sup>۱۰</sup> (دریا، جنگل، خیابان، کلاس درس و مواردی از این دست)

۲. دسته اجسام<sup>۱۱</sup> موجود در صحنه (صندلی، مرد، گربه و مواردی از این دست)

۳. ارتباط مکانی بین اجسام موجود (بالا، کنار، پشت و مواردی از این دست)

۴. رخدادی<sup>۱۲</sup> که در صحنه در حال اتفاق است (مانند نشستن، دویدن، کارکردن و مواردی از این دست)

<sup>۹</sup>Representation

<sup>۱۰</sup>Scence Class

<sup>۱۱</sup>Object Class

<sup>۱۲</sup>Event

### ۱-۳-۱ پژوهش‌های انجام‌شده در زمینه درک صحنه توسط مغز انسان

مساله درک صحنه، مانند بیشتر مسائل موجود در زمینه بینایی ماشین، الهام گرفته از نحوه رفتار انسان‌ها است. اغلب انسان‌ها با دیدن یک تصویر قادرند توصیف کامل و دقیقی از آن تصویر ارائه دهند که شامل تمام نکات لازم و ضروری نهفته در تصویر باشد. در بیشتر موارد، زمان مورد نیاز برای مغز انسان به منظور پردازش یک تصویر و توصیف آن، زمان بسیار کم و ناچیزی است. این مطلب، این ایده را در ذهن تداعی می‌کند که بخش قابل توجهی از اطلاعات مورد نیاز از هر تصویر، در اولین لحظاتی که تصویر به مغز می‌رسد (در نگاه اول) قابل استخراج است. بنابراین سامانه‌های رایانه‌ای باید قادر باشند با الگو گرفتن از مغز انسان، در کوتاه‌ترین زمان ممکن، اطلاعات کافی و مفید نهفته در تصویر را استخراج کرده و صحنه به نمایش کشیده شده در تصویر را توصیف کنند.

این فرض که مغز انسان می‌تواند در کوتاه‌ترین زمان ممکن، بیشترین حجم اطلاعات تصویر را به درستی استخراج نماید، توسط پژوهش‌گران متعددی مورد ارزیابی قرار گرفته است. از جمله اولین پژوهش‌هایی که به بررسی این فرض پرداخته‌اند می‌توان به [۱۰] و [۱۱] اشاره کرد. در این پژوهش‌ها، با نشان دادن تصاویر به صورت دنباله‌ای<sup>۱۳</sup> به مجموعه‌ای از افراد، از آن‌ها خواسته شده تا بهترین و دقیق‌ترین توصیفی را که می‌توانند، برای تصاویری که دیده‌اند، بازگو کنند. در این دو پژوهش نتیجه گرفته شده است که انسان می‌تواند یک تصویر معمولی را در بازه زمانی کمتر از ۲۰۰ میلی‌ثانیه، تشخیص داده و آن را توصیف کند. اگرچه این زمان برای تشخیص و توصیف یک تصویر کافیست، زمان مورد نیاز برای به‌خاطرسباری تصویر بسیار بیشتر از این مقدار است.

در پژوهش [۱] آزمایش دیگری انجام شده که از اهمیت بسیاری برخوردار است. در پژوهش‌های قبلی، افرادی که تصاویر را توصیف می‌کردند، درباره موضوع کلی تصاویر اطلاعاتی داشتند. اما در این آزمایش، تصاویر مختلفی از دنیای واقعی که محدود به شرایط خاصی نبوده‌اند، بدون ارائه پیش‌فرض درباره موضوع، به افراد نمایش داده شده و از آن‌ها خواسته شده که تصویر را به بهترین شکل توصیف کنند. آزمایشات در این پژوهش، در دو مرحله انجام شده‌اند.

۱. توسط یک رایانه، تصاویر متعددی در بازه‌های زمانی متفاوت به افراد نمایش داده می‌شوند و پس از اتمام زمان نمایش هر تصویر، یک ماسک بصری، تصویر را می‌پوشاند. در این حالت از افراد خواسته شده است که بهترین توصیف ممکن از تصویر را تایپ کنند. شرایط محیطی آزمایشات مطابق با استانداردها رعایت شده است. هر تصویر به طور تصادفی بین ۲۷ الی ۵۰۰ میلی ثانیه روی نمایش‌گر نمایش داده شده و سپس یک ماسک روی تصویر قرار گرفته و افراد فرصت دارند تا توصیف خود را از تصویر، بنویسند.

<sup>۱۳</sup>Image Series



PT = 107 ms

PT = 500 ms

This is outdoors. A black, furry dog is running/walking towards the right of the picture. His tail is in the air and his mouth is open. Either he had a ball in his mouth or he was chasing after a ball. (Subject EC)

I saw a black dog carrying a gray frisbee in the center of the photograph. The dog was walking near the ocean, with waves lapping up on the shore. It seemed to be a gray day out. (Subject JB)

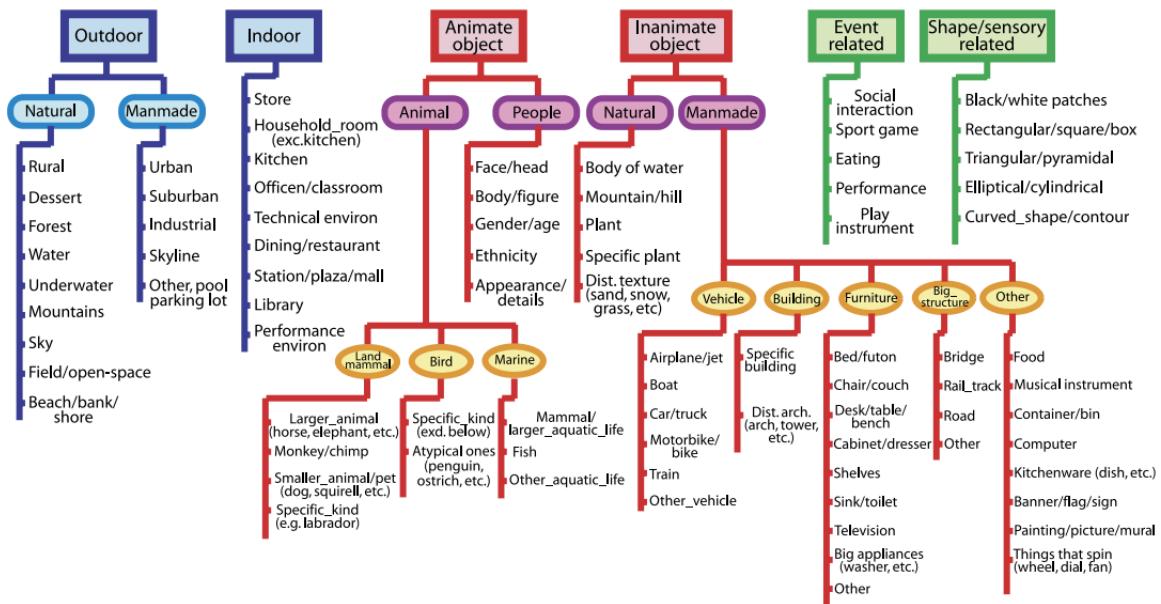


Inside a house, like a living room, with chairs and sofas and tables, no ppl. (Subject HS)

A room full of musical instruments. A piano in the foreground, a harp behind that, a guitar hanging on the wall (to the right). It looked like there was also a window behind the harp, and perhaps a bookcase on the left. (Subject RW)

شکل ۱: نمونه توصیف‌های افراد برای تصاویر [۱]

۲. در این مرحله، آزمایش روی افراد متفاوتی انجام شده است. این گروه افراد موظفند پس از دیدن تصاویر، به بهترین شکل ممکن آن‌ها را دسته‌بندی کنند. برخلاف افراد شرکت‌کننده در آزمایش قبلی که می‌توانستند به هر شکلی اطلاعات استخراج شده را بنویسند، به افراد حاضر در این گروه یک فرم مشخص از دسته‌اطلاعات مطلوب داده شده است که افراد موظفند آن را براساس محتوای تصویری که دیده‌اند، پر کنند. شکل ۲ ساختار ۲ ساختار مطلوب پاسخ افراد را در این آزمایش نمایش می‌دهد.



شکل ۲: ساختار مطلوب اطلاعات استخراج شده از تصاویر [۱]

این ساختار با تحلیل پاسخ‌های جمع‌آوری شده از آزمایش اول استخراج شده است و شامل انواع مختلفی از اطلاعات است که افراد در آزمایش اول به آن اشاره کرده‌اند.

شکل ۳ چند نمونه از تصاویر مورد استفاده در آزمایشات این پژوهش را نمایش می‌دهد. این تصاویر از اینترنت استخراج شده‌اند. برای استخراج این تصاویر از فضای اینترنت، از یک گروه افراد شامل ۱۰ نفر که با موضوع پژوهش

آشنا نبوده‌اند خواسته‌شده تا هر یک، نام ۵ دسته صحنه مختلف را به طور تصادفی بنویسند. پس از حذف نام‌های تکراری، ۳۰ نام منحصر به فرد باقی مانده‌است. سپس تصاویر مربوط به هریک از این نام‌ها توسط موتور جستجوی گوگل استخراج شده و ۶ تصویر از صفحات اولیه نتایج به عنوان تصاویر نمونه انتخاب شده‌اند.



(آ) چند نمونه از تصاویر در محیط باز  
(ب) چند نمونه از تصاویر در محیط بسته

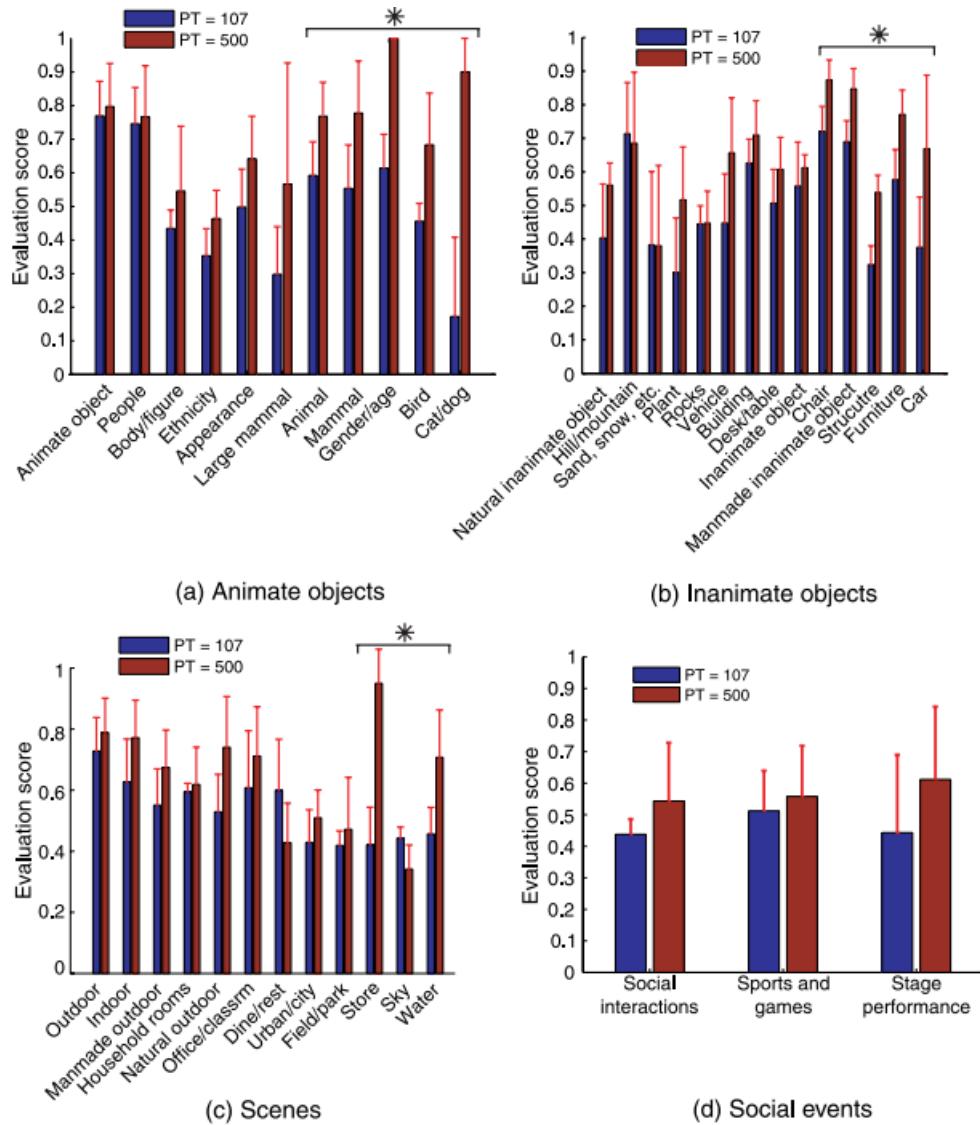
شکل ۲: تصاویر دنیای واقعی مورد استفاده در آزمایشات [۱]

ارزشمندترین نکته درباره پژوهش انجام‌شده، یافته‌های آن است. این پژوهش نکاتی را در مورد توانایی مغز انسان در توصیف صحنه روشن می‌کند که حائز اهمیت هستند. در ادامه این نتایج را بررسی خواهیم کرد.

### ۱-۳-۲ نتایج به دست آمده از آزمایشات

۱. حداقل زمان لازم برای مغز انسان به منظور درک صحنه، برابر با ۵۰۰ میلی ثانیه است.
۲. این مدت زمان، برای صحنه‌های ساده و بدون پیچیدگی، به حدود ۱۰۰ میلی ثانیه می‌رسد. به عنوان نمونه در شکل ۱ تصویر اول که دارای پیچیدگی‌های کمتری نسبت به تصویر دوم است در مدت زمان ۱۰۷ میلی ثانیه، به‌طور کامل توصیف شده‌است در صورتی که تصویر دوم که به نسبت، پیچیده‌تر است، مدت زمان بیشتری برای توصیف نیاز داشته‌است.
۳. با استفاده از ساختارمندسازی پاسخ‌های افراد در آزمایش دوم و اطلاعات جمع‌آوری شده در درخت پاسخ‌ها (که در شکل ۲ نمایش داده شده است) و میانگین‌گیری روی تمام تصاویر، نمودارهای مقایسه‌ای برای مدت

زمان ۱۰۷ میلی ثانیه و ۵۰۰ میلی ثانیه ایجاد شده است. شکل ۴ نمودارهای مقایسه‌ای را نمایش می‌دهد. در این نمودارها، میله‌های قرمز نشان‌دهنده نتایج برای زمان ۵۰۰ میلی ثانیه و میله‌های آبی نمایش‌دهنده نتایج برای حالت ۱۰۷ میلی ثانیه هستند. در دو نمودار اول (نمودارهای بالا سمت راست و بالا سمت چپ) تشخیص و استخراج اطلاعات مربوط به اجسام مختلف بسته به متحرک بودن<sup>۱۴</sup> یا متحرک نبودن<sup>۱۵</sup> آن‌ها، در نمودار سوم (نمودار پایین سمت چپ) تشخیص و استخراج اطلاعات مربوط به صحنه موجود در تصویر و در نمودار چهارم (نمودار پایین سمت راست) تشخیص و استخراج اطلاعات مربوط به رخداد موجود در تصویر، مورد بررسی قرار گرفته‌اند.



شکل ۴: نمودارهای مقایسه‌ای عملکرد مغز انسان در درک صحنه در بازه‌های زمانی ۱۰۷ و ۵۰۰ میلی ثانیه<sup>[۱]</sup>

همان‌طور که مشخص است، مدت زمان ۱۰۷ میلی ثانیه برای مغز انسان، زمان بهینه برای توصیف صحنه است. تفاوت‌های بین نتایج در اکثر موارد، جزئی و در مقابل تفاوت زمانی موجود، بسیار کوچک هستند. به علاوه، در تمام مواردی که نیاز به اطلاعات کلی از تصویر وجود دارد، تفاوت بین دو بازه زمانی چندان

<sup>۱۴</sup> Animated

<sup>۱۵</sup> Inanimate

چشم‌گیر نیست، اما در مواردی که برای تشخیص نیاز به دانستن جزئیات بیشتر از تصویر وجود دارد (مانند سن، جنسیت و نوع حیوان) تفاوت بین دو زمان، قابل ملاحظه است.

همین‌طور با مقایسه تفاوت عملکرد بین حالات متحرك بودن و متتحرك نبودن اجسام، فواصل موجود در نمودارها قابل ملاحظه می‌شود. در حالت کلی، تفاوت بین عملکرد مغز در دو بازه، در حالتی که اجسام ساکن در تصویر وجود دارند به مراتب کمتر از حالتی است که اجسام موجود در تصویر، متتحرك باشند.

شکل ۵ نمونه دیگری از نتایج به دست آمده از آزمایشات را در مدت‌زمان‌های مختلف نمایش می‌دهد.

		
PT 27 ms	There was a range of dark splotches in the middle of the picture, running from most of the way on the left side, to all the way on the right side. This was surrounded primarily by a white or light gray color. (Subject: KM)	Couldn't see much; it was mostly dark w/ some square things, maybe furniture. (Subject: AM)
PT 40 ms	I saw a very bright object, shaped in a pyramidal shape. There was something black in the front, but I couldn't tell what it was. (Subject: JB)	Looked like something black in the center with four straight lines coming out of it against a white background. (Subject: AM)
PT 67 ms	Possibly outdoors. maybe a few ducks, or geese. Water in the background. (Subject: JL)	This looked like an indoor shot. Saw what looked like a large framed object (a painting?) on a white background (i.e., the wall). (Subject: RW)
PT 500 ms	It was definitely on a coast by the ocean with a large [r]ock in the foreground and at least three birds sitting on the rock. (Subject: CC)	The first thing I could recognize was a dark splotch in the middle. It may have been rectangular-shaped, with a curved top...but, that's just a guess. (Subject: KM)
	I saw the interior of a room in a house. There was a picture to the right, that was black, and possibly a table in the center. It seemed like a formal dining room. (Subject: JB)	A person, I think, sitting down or crouching. Facing the left side of the picture. We see their profile mostly. They were at a table or were some object was in front of them (to their left side in the picture). (Subject: EC)
	Some fancy 1800s living room with ornate single seaters and some portraits on the wall. (Subject: WC)	This looks like a father or somebody helping a little boy. The man had something in his hands, like a LCD screen or laptop. they looked like they were standing in a cubicle. (Subject: WC)

شکل ۵: نمونه‌ای از نتایج به دست آمده از آزمایشات [۱]

## ۴-۱ جمع‌بندی

با توجه به افزایش چشم‌گیر تعداد تصاویر مورد استفاده کاربران در فضاهای مجازی و همین‌طور با در نظر گرفتن گرایش روزافزون کاربران به ذخیره‌سازی تصاویر در رایانه‌های شخصی، مساله مدیریت این تصاویر و یافتن تصاویر خاص بین مجموعه تصاویر موجود، به یکی از مسائل مهم و پرکاربرد در زمینه بینایی ماشین تبدیل شده است. گام اساسی در این راستا، دست‌یابی به سامانه‌ای است که قادر به تولید خودکار شرح برای تصاویر باشد. شرح این تصاویر که در قالب جملات زبان طبیعی ارائه می‌شود باید علاوه بر سازگاری با موضوع تصویر و توصیف صحیح صحنه، به لحاظ دستور زبان و معنا صحیح و کامل باشد. فرایند تولید خودکار شرح برای تصاویر، از دو مرحله اصلی تشکیل می‌شود:

۱. نگاشت تصویر ورودی به فضای بردار ویژگی‌ها (درک صحنه)

۲. تولید جملات زبان طبیعی مبتنی بر محتوای بردار ویژگی‌ها

مساله درک صحنه، یکی از چالش برانگیزترین مسائل در زمینه بینایی ماشین است. با این وجود، تا کنون تعریف دقیق و کاملی از این مفهوم ارائه نشده است. به طور کلی می‌توان درک صحنه را فرایندی دانست که طی آن اطلاعات بصری موجود در تصویر استخراج شده و در قالب خاصی بازنمایی می‌شوند. میزان و نوع این اطلاعات را نمی‌توان به طور کلی تعریف کرد. حوزه تعریف اطلاعات و کیفیت مطلوب آن‌ها بسته به کاربرد در هر حوزه تعریف می‌شود.

در بین پژوهش‌های مربوط به تولید خودکار شرح برای تصاویر، انواع اطلاعات مطلوب، عموماً شامل موارد زیر می‌شود:

۱. دسته صحنه

۲. دسته اجسام

۳. ارتباط مکانی بین اجسام موجود

۴. رخدادی که در صحنه در حال اتفاق است

پژوهش‌گران از گذشته بر این عقیده بوده‌اند که مغز انسان در اولین لحظات مشاهده یک تصویر، قادر است اطلاعات کافی و مفید برای درک صحنه را استخراج کند. پژوهش‌های متعددی در این زمینه انجام شده‌اند که هریک به بررسی جوانب خاصی از این فرضیه پرداخته‌اند. به عنوان نمونه، پژوهش [۱۰] و [۱۱] با استفاده از دنباله‌های تصاویر، مدت زمان مورد نیاز برای مغز انسان به جهت درک صحنه را کمتر از ۲۰۰ میلی‌ثانیه تخمین زده‌اند.

در پژوهش [۱]، یک آزمایش دو مرحله‌ای برای بررسی تاثیر مدت زمان مشاهده تصاویر بر عملکرد مغز در توصیف صحنه، انجام شده است. در این آزمایش که در دو مرحله انجام شده، ابتدا گروهی از افراد با دیدن تصاویر در مدت زمان بین ۲۷ تا ۵۰۰ میلی‌ثانیه، موظف به توصیف تصویر بوده‌اند. سپس گروه دیگری از افراد با دیدن تصاویر در مدت زمان‌های مختلف، ملزم به پر کردن فرم از پیش تعیین‌شده‌ای بودند که با توجه به پاسخ‌های بهدست‌آمده از آزمایش اول، تدوین شده است.

نتایج این آزمایشات نشان می‌دهد، مدت زمان ۱۰۷ میلی‌ثانیه برای تشخیص و بخش قابل توجهی از اطلاعات موجود در تصویر کافیست؛ اگرچه، در مواردی که دق به جزئیات ضروری است (مانند تشخیص سن، جنسیت، نوع حیوان) و برای تشخیص و استخراج اطلاعات اجسام متحرک، مدت زمان ۵۰۰ میلی‌ثانیه، بهبود قابل توجهی در عملکرد مغز ایجاد می‌کند.

## ۲ فصل دوم

### درگ صحنه

## ۱-۲ درک صحنه

درک صحنه یکی از چالش‌های اساسی در زمینه بینایی ماشین است که روش‌های مختلفی برای دست‌یابی به آن ارائه شده است. با وجود تعدد پژوهش‌های موجود در این مورد، ارائه تعریف جامع و شامل برای این مفهوم کاری بسیار دشوار است. عموماً این مفهوم، بسته به مورد کاربرد و هدف پژوهش، به استخراج مجموعه مشخصی از اطلاعات در مورد صحنه که برای پژوهش، کافی و مفید باشد محدود می‌شود. به همین دلیل، مجموعه اطلاعات مطلوب از تصویر که باید استخراج شود در هر پژوهش به طور خاص تعریف می‌شود.

درک صحنه در زمینه تولید خودکار شرح بر تصاویر، به طور عام شامل موارد زیر می‌شود:

۱. تشخیص اجسام موجود در صحنه و دسته‌بندی آن‌ها (مانند توپ، تلویزیون)

۲. تشخیص ارتباط مکانی بین اجسام موجود در صحنه (مانند پشت، بالا)

۳. دسته‌بندی محیط (مانند جنگل، دریا)

۴. دسته‌بندی فعالیت به تصویر کشیده شده (مانند راه‌رفتن، خوابیدن)

## ۲-۲ روش‌های مختلف موجود

فعالیت‌های متعددی برای تشخیص هر یک از موارد بالا انجام شده است. به طور عام می‌توان روش‌های مورد استفاده در استخراج اطلاعات مطلوب صحنه را در زمینه تولید خودکار شرح بر تصاویر به دو دسته عمدی زیر تقسیم‌بندی نمود:

۱. استفاده از مدل‌های گرافی احتمالی<sup>۱۶</sup>

در این دسته از روش‌ها، با استفاده از مدل‌های گرافی احتمالی در مورد حضور یا عدم حضور اجسام مختلف در صحنه و رابطه بین اجسام موجود استنتاج نمود. همین‌طور فرایندهایی مانند قطعه‌بندی تصویر<sup>۱۷</sup> در این روش‌ها با استفاده از مدل‌های گرافی احتمالی انجام می‌شوند. به عنوان نمونه، در مقاله [۱۲] یک مدل

<sup>۱۶</sup>Probabilistic Graphical Models (PGMs)

<sup>۱۷</sup>Image Segmentation

میدان تصادفی شرطی<sup>۱۸</sup> برای تجزیه معنایی<sup>۱۹</sup> تصویر ارائه شده است که با استفاده از آن می‌توان در مورد حضور یا عدم حضور اجسام مختلف به طور توان در صحنه تصمیمگیری کرد.

۲. استفاده از شبکه‌های عصبی کانولوشنی عمیق در این دسته از روش‌ها، با استفاده از شبکه‌های عصبی کانولوشنی عمیق، پس از قطعه‌بندی تصاویر، اقدام به تفکیک اجسام مختلف در صحنه و برچسب‌گذاری هر جسم، بسته به یادگیری انجام شده، می‌شود. به عنوان نمونه در مقاله [۷] یک شبکه عصبی کانولوشنی عمیق معرفی شده است که قادر به برچسب‌گذاری اجسام مختلف در صحنه است. برچسب‌های مورد استفاده در این پژوهش، عبارات مختلف موجود در جملات توصیف‌گر هر تصویر در مجموعه‌دادگان هستند.

نمونه‌های متعددی از این دست پژوهش‌ها، در هر دسته، انجام شده است که در ادامه چند مورد از آن‌ها بررسی خواهد شد.

### ۳-۲ روش‌های مبتنی بر مدل‌های گرافی احتمالی

همان‌طور که قبلاً ذکر شد، روش‌های مبتنی بر استفاده از مدل‌های گرافی احتمالی، از جمله پرکاربردترین روش‌ها در مرحله درک صحنه در زمینه تولید خودکار شرح بر تصاویر هستند. این روش‌ها با استفاده از نظریه گراف، آمار و احتمالات اقدام به ارائه یک توزیع احتمالی برای پارامتر مورد بررسی، با توجه به داده‌های موجود در مجموعه آموزشی می‌کنند. مدل‌های استاندارد مختلفی در پژوهش‌ها مورد استفاده قرار می‌گیرند که تعدادی از آن‌ها به عنوان نمونه در این بخش مورد بررسی قرار خواهند گرفت.

#### ۳-۲-۱ استفاده از مدل میدان تصادفی مارکف<sup>۲۰</sup>[۱۳]

مقاله [۱۳] با استفاده از یک مدل ساده میدان تصادفی مارکف، فرایند درک صحنه را انجام می‌دهد و با استفاده از همین مدل، اقدام به تولید جملات توصیف‌گر تصویر می‌نماید. در این فصل به بررسی فرایند درک صحنه در این مقاله می‌پردازیم و بررسی فرایند تولید جمله را به فصل بعدی موکول می‌نماییم.

درک صحنه در این پژوهش محدود به ارتباط بین سه مفهوم در هر تصویر شده است؛ به این معنی که به ازای هر تصویر، یک سه‌تایی «جسم، فعالیت، صحنه»<sup>۲۱</sup> ایجاد می‌شود که بیان‌کننده اطلاعات مطلوب موجود در تصویر است. میدان<sup>۲۲</sup> «جسم»، دربر دارنده برچسب حاصل از دسته‌بندی اجسام موجود در صحنه، میدان «فعالیت»، دربر دارنده اطلاعات مربوط به فعالیت در حال انجام و میدان «صحنه» دربردارنده اطلاعات مربوط به محیط تصویر هستند. به فضای سه‌تایی‌های ایجاد شده برای اطلاعات مطلوب در درک صحنه، فضای معنا<sup>۲۳</sup> می‌گویند.

شکل ۶ نمایی از نگاشت اطلاعات از فضای تصاویر و جملات به فضای معنایی، نمایش می‌دهد. همان‌طور که در شکل مشخص است، به ازای هر تصویر، یک سه‌تایی معنایی ایجاد می‌شود. همین‌طور به ازای هر جمله در

<sup>۱۸</sup>Conditional Random Field (CRF)

<sup>۱۹</sup>Semantic Parsin g

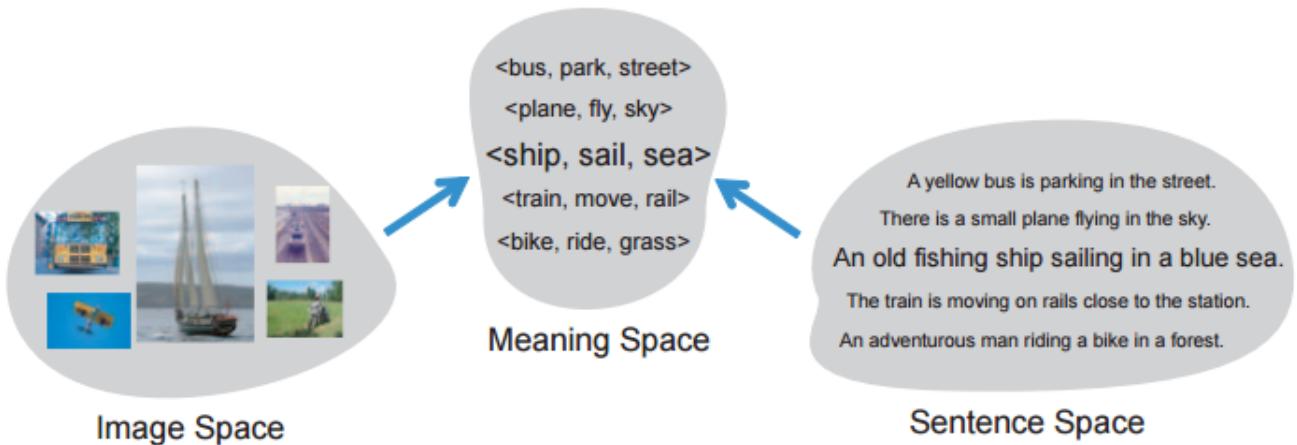
<sup>۲۰</sup>Markov Random Field (MRF)

<sup>۲۱</sup><Object, Activity, Scene>

<sup>۲۲</sup>Field

<sup>۲۳</sup>Meaning Space

فضای جملات، یک سه‌تایی ایجاد می‌شود به‌طوری که جملات و تصاویر متناظر شان، به یک سه‌تایی یکسان، نگاشت شوند. همان‌طور که مشخص است، با داشتن نگاشتهایی که خواص مذکور را داشته باشند، می‌توان با استفاده از سه‌تایی‌های فضای معنا، تصاویر را مدیریت کرد.



شکل ۶: نگاشت تصویر به فضای معنایی. فضای معنایی شامل اطلاعات مطلوب برای استخراج در فرایند درک صحنه است. به ازای هر تصویر، یک سه‌تایی ایجاد می‌شود [۱۳].

مدل میدان تصادفی مارکف مورد استفاده در این پژوهش، یک مدل کوچک و ساده، شامل ۳ گره است. شکل ۷ طرح‌واره‌ای از مدل میدان تصادفی مارکف مورد استفاده در این پژوهش را نمایش می‌دهد. همان‌طور که در شکل مشخص است، به ازای هر کدام از میدان‌های تعریف شده در فضای معنایی، یک گره در این مدل وجود دارد. مقادیر مختلف در هر گره، برابر است با مقادیر مختلف موجود در میدان متناظر، در فضای معنا که با توجه به داده‌های مجموعه آموزشی مشخص می‌شوند. همین‌طور به ازای هر دو گره موجود در این مدل، یک یال بیان‌کننده ارتباط بین دو میدان در فضای معنایی وجود دارد.

برای استنتاج در این مدل، لازم است ابتدا فاکتورهای مورد استفاده در مدل را شناخته و مقادیر آن‌ها را مشخص نماییم. در مدل پیشنهادی، دو نوع فاکتور تعریف شده است:

#### ۱. فاکتورهای گره

این فاکتورها، برای مشخص کردن میزان شباهت مقادیر مختلف گره با تصویر ورودی، تعریف شده‌اند. ویژگی‌های مورد استفاده برای مقداردهی این فاکتورها، شامل موارد زیر هستند:

(آ) استفاده از آشکارکننده‌های<sup>۲۴</sup> فلزنسوالب<sup>۲۵</sup>، به منظور محاسبه امتیاز اطمینان<sup>۲۶</sup> برای هر دسته از اجسام موجود در مجموعه داده [۱۴].

پس از محاسبه امتیاز اطمینان همه دسته‌های موجود، دسته‌ای که بیشترین امتیاز را دارد می‌تواند به عنوان دسته منتخب در میدان متناظر گره، انتخاب شود. در فرایند مقداردهی این ویژگی، قبل از انجام محاسبات، اطمینان حاصل می‌شود که از هر دسته موجود، حداقل یک تصویر در مجموعه داده وجود داشته باشد.

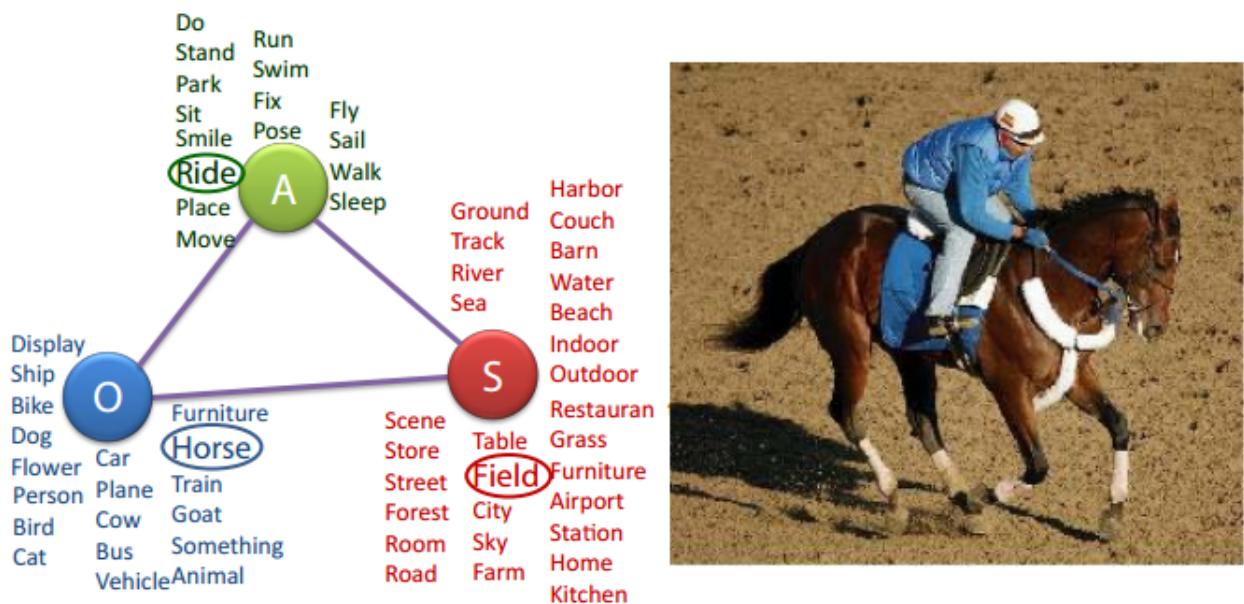
<sup>۲۴</sup>Detector

<sup>۲۵</sup>Felzenszwaalb

<sup>۲۶</sup>Confidence Score

(ب) استفاده از پاسخ دسته‌بندی‌کننده دیوالا<sup>۷۷</sup>، ارائه شده در مقاله [۱۵]

(ج) استفاده از دسته‌بندی‌کننده مبتنی بر گیست<sup>۷۸</sup>



شکل ۷: طرحواره مدل میدان تصادفی مارکف ارائه شده در پژوهش [۱۲] که شامل ۳ گره است. در این مدل، به ازای هر میدان از فضای معنا، یک گره وجود دارد و بین هر سه گره، به طور دو به دو، یک یال موجود است [۱۳].

بر اساس مقادیر محاسبه شده برای ویژگی‌های بالا و با استفاده از الگوریتم ماشین بردار پشتیبان<sup>۷۹</sup>، یک دسته‌بندی برای هر گره ارائه می‌شود که بیان‌کننده دسته‌ویژگی‌های مربوط به مقادیر مختلف گره است. با استفاده از این دسته‌بندی، با ورود هر تصویر، می‌توان برای هر مقدار در هر گره، یک امتیاز شباهت محاسبه نمود. استفاده از الگوریتم یافتن نزدیک‌ترین همسایه‌های موجود برای هر تصویر ورودی، بر اساس امتیاز شباهت محاسبه شده و میانگین‌گیری روی همسایه‌های استخراج شده، معیار خوبی از تخمین مقدار هر گره، به ازای هر تصویر ورودی ایجاد می‌کند. به این ترتیب، با ورود هر تصویر می‌توان برای هر کدام از گره‌های موجود در مدل، یک مقدار محتمل مشخص نمود. سه‌تایی شامل مقادیر محتمل بدست‌آمده در هر گره، سه‌تایی متناظر تصویر ورودی در فضای معنا را مشخص می‌کند.

## ۲. فاکتور یال

این فاکتور، برای مشخص کردن میزان ارتباط مقادیر مختلف دو گره با یکدیگر در تصویر ورودی مورد استفاده قرار می‌گیرند.

<sup>۷۷</sup>divvala

<sup>۷۸</sup>Gist-based classification response

<sup>۷۹</sup>Support Vector Machine (SVM)

## ۲-۳-۲ استفاده از مدل میدان تصادفی شرطی<sup>۳۰</sup>

در این پژوهش، مساله درک صحنه در قالب یک مساله استنتاج با استفاده از مدل میدان تصادفی شرطی بیان شده است. مدل میدان تصادفی شرطی، یکی از پرکاربردترین مدل‌های گرافی احتمالی در زمینه درک صحنه است که پژوهش‌های متعددی از آن به عنوان مدل اصلی در درک صحنه استفاده کرده‌اند. به عنوان نمونه، در مقاله‌های [۱۶] و [۱۷] از مدل میدان تصادفی شرطی به منظور توصیف صحنه استفاده شده است.

پژوهش [۱۶] سعی در توصیف اجسام سه‌بعدی با استفاده از قطعه‌بندی تصاویر دو بعدی، هندسه سه‌بعدی و روابط بین صحنه و اجسام موجود، دارد. در این پژوهش، پس از استخراج ویژگی‌ها و اطلاعات بدست‌آمده از منابع مختلف، عمل استنتاج توسط یک مدل تصادفی شرطی انجام می‌شود که منجر به نگاشت تصویر ورودی به فضای معنایی می‌شود. همین‌طور در پژوهش [۱۷]، یک چارچوب کاری<sup>۳۱</sup> احتمالی برای استنتاج درباره نواحی مختلف تصویر، اجسام موجود و ویژگی‌های مختلف آن‌ها مانند دسته‌بندی، موقعیت مکانی و ابعاد، مبتنی بر مدل میدان تصادفی شرطی، ارائه شده است. با توجه به وسعت و تعدد فعالیت‌های انجام شده، در این بخش، مرحله درک صحنه یک پژوهش انجام شده در زمینه تولید خودکار شرح بر تصاویر را مورد بررسی قرار می‌دهیم. لازم به ذکر است، مرحله تولید جملات توصیف‌کننده پژوهش مورد بحث، در فصل تولید جملات زبان طبیعی مورد بررسی قرار خواهد گرفت.

در پژوهش [۱۲] از مدل میدان تصادفی شرطی برای توصیف صحنه و اجسام موجود در آن استفاده شده است. میدان‌های تصادفی در این مدل، شامل متغیرهای زیر هستند:

۱. متغیرهای تصادفی بیان‌کننده برچسب دسته متناظر قطعات مختلف هر تصویر به شیوه سلسله مراتبی دارای دو سطح

۲. متغیرهای تصادفی باینری بیان‌کننده صحت دسته تشخیص داده شده برای هر جسم

شکل ۸ طرح‌واره مدل سلسله‌مراتبی ارائه شده در پژوهش [۱۲] را نمایش می‌دهد. همان‌طور که مشاهده می‌شود این مدل از دو سطح انتزاع، یکی برای برچسب قطعات مختلف تصویر و دیگری برای حضور یا عدم حضور هر دسته از اجسام در تصویر، تشکیل شده است.

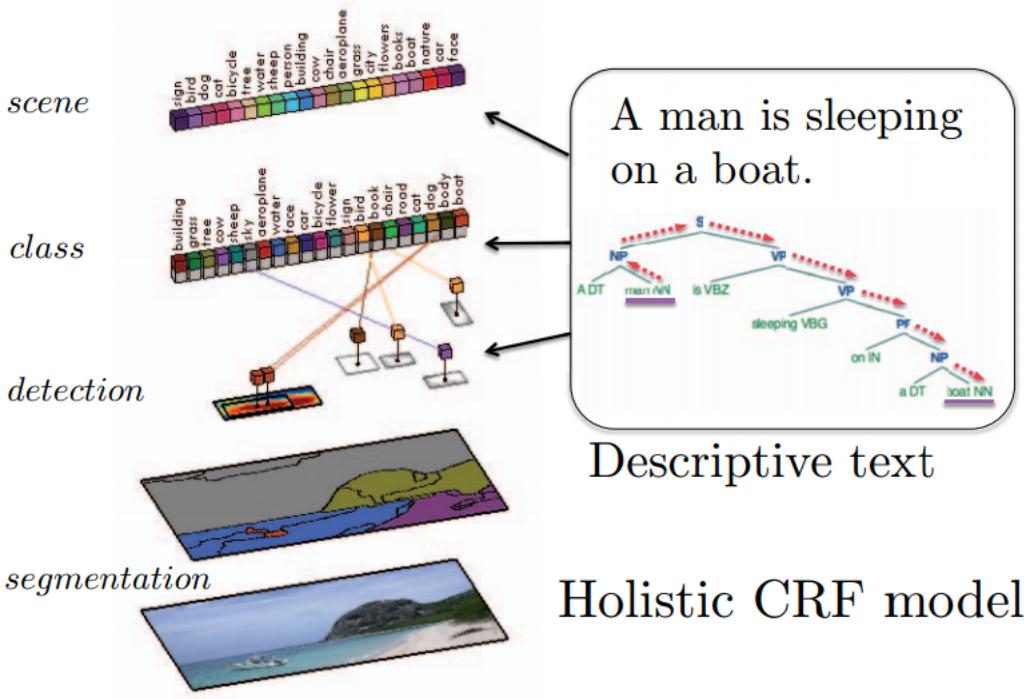
دو دسته متغیر تصادفی مختلف، که هر یک نماینده متغیرهای تصادفی موجود در یکی از این سطوح انتزاع هستند، تعریف شده‌اند؛ متغیرهای تصادفی  $C$ ،  $X_i \in \{1, \dots\}$  بیان‌کننده دسته قطعه  $i$  از سطح پایین سلسله مراتب و متغیرهای تصادفی  $C$ ،  $Y_j \in \{1, \dots\}$  بیان‌کننده دسته قطعه  $j$  از سطح بالای سلسله مراتب. به علاوه، دو دسته متغیر تصادفی دیگر به نام‌های  $b_k$  و  $z_k$  به ترتیب برای نمایش حضور یا عدم حضور یک تشخیص کاندید<sup>۳۲</sup> و حضور یا عدم حضور جسم با دسته  $k$  در تصویر، تعریف شده‌اند. با توجه به متغیرهای تعریف شده، مدل کلی میدان تصادفی شرطی را می‌توان معادل رابطه<sup>۳۳</sup> ۱ تعریف کرد. در این رابطه  $(a_\alpha)^{\Psi_\alpha^{type}}$  نماینده تابع پتانسیل تعریف شده روی متغیرهای مختلف است. با این تعریف، یافتن تخمین MAP<sup>۳۴</sup>، منجر به یافتن پاسخ مورد نظر می‌شود.

<sup>۳۰</sup> Conditional Random Field (CRF)

<sup>۳۱</sup> Framework

<sup>۳۲</sup> Candidate Detection

<sup>۳۳</sup> MAP Estimation



شکل ۸: طرح‌واره مدل سلسله مراتبی مبتنی بر میدان تصادفی شرطی که بر اساس اطلاعات بصری و اطلاعات جملات توصیف‌کننده شرح محتمل تصویر را تولید می‌نماید [۱۲].

در ادامه، توابع پتانسیل مختلف که در این پژوهش تعریف شده‌اند، ارائه خواهد شد. لازم به ذکر است در تمام این موارد، برای سهولت، توابع پتانسیل به شکل لگاریتمی تعریف شده‌اند.

$$P(X, Y, b, z) = \frac{1}{Z} \prod_{type} \prod_{\alpha} \Psi_{\alpha}^{type}(a_{\alpha}) \quad (1)$$

توابع پتانسیل مختلف تعریف شده در این پژوهش عبارتند از:

#### ۱. پتانسیل قطعه‌بندی یگانی <sup>۳۴</sup>

پتانسیل قطعه‌بندی یگانی در هر قطعه و هر ابرقطعه <sup>۳۵</sup> از تصویر، با استفاده از میانگین‌گیری روی امتیاز افزایش تکستون <sup>۳۶</sup> که در پژوهش [۱۸] ارائه شده است، انجام می‌شود.

#### ۲. انطباق بین متغیرهای دو سطح انتزاع با یکدیگر

یک مقدار جریمه به ازای دسته‌های مخالف بین دو سطح در نظر گرفته می‌شود تا در حد امکان، دسته‌های منتخب از بین سطوح مختلف، با یکدیگر انطباق داشته باشند. پتانسیل تعریف شده در این بخش معادل رابطه ۲ تعریف می‌شود.

$$\phi_{ij}(X_i, Y_j) = \begin{cases} -\gamma & X_i \neq Y_j \\ 0 & X_i = Y_j \end{cases} \quad (2)$$

در رابطه ۲، پارامتر  $\gamma$  در فرآیند یادگیری که منجر به بهینه‌سازی پارامترهای مختلف مدل می‌شود، به دست می‌آید.

<sup>۳۴</sup>Unary Segmentation Potential

<sup>۳۵</sup>Supersegment

<sup>۳۶</sup>Texton Boost

### ۳. پتانسیل انطباق تصویر و دسته جسم

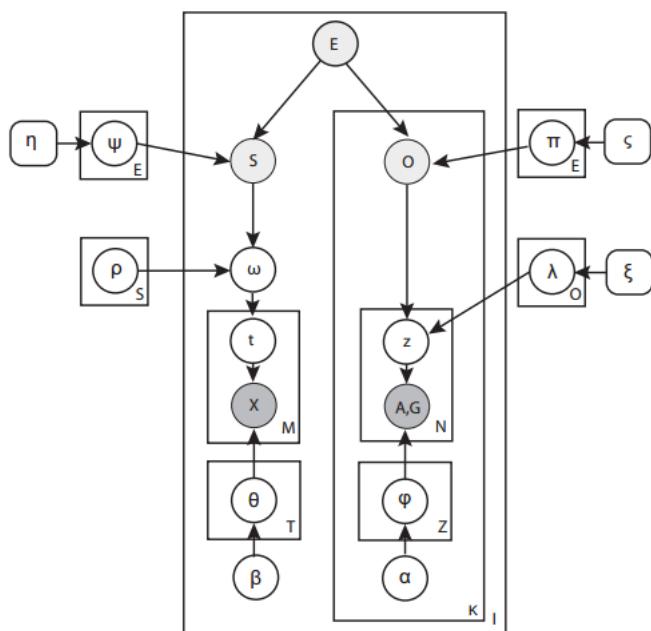
برای اندازه‌گیری میزان انطباق هر کدام از دسته‌های موجود برای اجسام با تصویر ورودی، از معیار انطباق ارائه شده در پژوهش [۱۹] توسط فلزنسوالب که به روش دی پی ام <sup>۳۷</sup> مشهور است، استفاده شده است. برای کاهش تعداد پارامترها و افزایش کارایی مدل استفاده شده، برای هر تصویر حداقل ۳ دسته جسم، به عنوان دسته‌های منتخب کاندید، در نظر گرفته می‌شوند.

### ۳-۳-۲ استفاده از سایر مدل‌های گرافی احتمالی

در بین پژوهش‌های موجود در زمینه درک صحنه با استفاده از روش‌های احتمالاتی، علاوه بر مدل‌های استاندارد، از مدل‌های مولد دیگر در پژوهش‌های متعددی استفاده شده است. در ادامه این بخش، به بررسی چند نمونه از این مدل‌ها خواهیم پرداخت.

#### ۱. دسته‌بندی تصاویر بر اساس صحنه و اجسام موجود به طور توأم [۲]

مدل استفاده شده در این پژوهش، از تصاویر در سطح صحنه و سطح اجسام استفاده کرده و با یکپارچه‌سازی و تجمعی اطلاعات موجود در این دو سطح، اقدام به دسته‌بندی تصویر می‌نماید. شکل ۹ مدل استفاده شده در این پژوهش را به منظور یکپارچه‌سازی و تجمعی اطلاعات حاصل از تحلیل صحنه و تشخیص اجسام موجود در آن، ارائه می‌دهد.



شکل ۹: مدل استفاده شده به منظور تجمعی اطلاعات صحنه و اجسام موجود در آن به منظور دسته‌بندی تصاویر [۲]

یکی از اهدافی که در این پژوهش دنبال می‌شود، برحسب گذاری معنایی<sup>۳۸</sup> تمام پیکسل‌های موجود در تصویر است. به همین منظور، تمام تصاویر مورد استفاده، به نواحی  $10 * 10$  تقسیم شده و مورد استفاده قرار

<sup>۳۷</sup>DPM

<sup>۳۸</sup>Semantic Labelling

می‌گیرند. برای بررسی بهتر مدل، ابتدا متغیرهای تصادفی مورد استفاده را تعریف کرده و سپس به بررسی روند یادگیری و استنتاج مدل می‌پردازیم.

متغیر تصادفی  $X$  که حاوی اطلاعاتی مبتنی بر حضور یا عدم حضور دسته‌های مختلف صحنه است، در بخش تشخیص صحنه به کار می‌رود. اطلاعات این متغیر با استفاده از توصیف‌کننده سیفت<sup>۳۹</sup> و به ازای هر ناحیه از تصویر، به دست می‌آید. برای بخش تشخیص اجسام موجود در صحنه، از دو منبع اطلاعاتی مختلف استفاده می‌شود. اطلاعات مربوط به حضور یا عدم حضور دسته‌های مختلف اجسام در متغیر تصادفی  $A$  و اطلاعات مربوط به شکل کلی آن‌ها در متغیر تصادفی  $G$  نمایش داده می‌شود.

هر گره از مدل ارائه شده، نماینده یک متغیر تصادفی است. گره‌هایی که با رنگ تیره مشخص شده‌اند، نماینده متغیرهایی هستند که در فرایند آموزش دیده می‌شوند و بقیه متغیرها، متغیرهای مخفی<sup>۴۰</sup> هستند. گره‌های خاکستری روشن‌تر، متغیرهایی هستند که فقط در فرایند آموزش دیده می‌شوند در حالی که متغیرهای تیره‌تر در هر دو فرایند آموزش و آزمون مشاهده می‌شوند.

متغیر تصادفی  $E$ ، نماینده یک دسته از رخداد<sup>۴۱</sup> های ممکن است. توزیع احتمال اولیه این متغیر تصادفی، یک توزیع یکنواخت فرض شده است که به هر تصویر ورودی، بر اساس همین توزیع، یک مقدار خاص از این متغیر تصادفی اختصاص داده می‌شود. با داشتن دسته رخداد موجود در تصویر، یک تصویر صحنه<sup>۴۲</sup> متناظر با تصویر ورودی تولید می‌شود. با فرض وجود  $S$  دسته صحنه مختلف در مجموعه‌داده، به هر تصویر، تنها یک دسته صحنه اختصاص داده می‌شود. روند اختصاص دسته صحنه به تصویر مطابق زیر است:

\* ابتدا یک دسته اولیه مطابق با توزیع احتمال شرطی  $P(S|E, \psi)$  به تصویر اختصاص داده می‌شود.  
 یک پارامتر چندجمله‌ای<sup>۴۳</sup> حاکم بر توزیع احتمالاتی  $S$  به شرط داشتن  $E$  است. به علاوه،  $\psi$  یک ماتریس به ابعاد  $S * E$  و پارامتر  $\theta$  یک بردار  $S$  بعدی در نقش مقدار اولیه دیریکله<sup>۴۴</sup> برای پارامتر  $\psi$  است.

\* در قدم بعدی با داشتن مقدار  $S$ ، پارامترهای  $\omega$  را بر اساس احتمال  $P(\omega|S, \rho)$  تولید می‌کنیم. از آن جا که  $\omega$  پارامتر چندجمله‌ای گره‌های مخفی  $t$  هستند، باید مجموع همه آن‌ها برابر با یک باشد. به علاوه،  $\rho$  یک ماتریس به ابعاد  $S * T$  و مقدار اولیه دیریکله برای پارامتر  $\omega$  است که در آن  $T$  تعداد کل  $t$ ‌ها است.

\* برای تولید هر یک از  $M$  ناحیه تصویر (مقادیر متغیر تصادفی  $X$ ) به شکل زیر عملی می‌کنیم:  
 - یک مقدار  $t$  از توزیع احتمال  $Mult(\omega)$  تولید می‌شود که مشخص‌کننده موضوعی<sup>۴۵</sup> است که این ناحیه از تصویر مطابق با آن تولید شده است.

- متغیر تصادفی  $X$  از توزیع احتمالی  $P(X|t, \theta)$  تولید می‌شود.  $\theta$  یک ماتریس به ابعاد  $V_s$

<sup>۳۹</sup>SIFT Descriptor

<sup>۴۰</sup>Latent

<sup>۴۱</sup>Event

<sup>۴۲</sup>Scene Image

<sup>۴۳</sup>Multinomial

<sup>۴۴</sup>Dirichlet prior

<sup>۴۵</sup>Topic

است که در آن  $V_s$  تعداد کلمات موجود در پایگاه داده مربوط به صحنه  $s$  است. به علاوه،  $\theta$  یک پارامتر چندجمله‌ای برای  $X$  است و  $\beta$  مقدار اولیه دیریکله برای  $\theta$ .

همانند فرایندی که طی آن، تصویر صحنه به تصویر ورودی اختصاص داده می‌شود، فرایندی وجود دارد که طی آن تصویر اجسام<sup>۴۶</sup> به تصویر ورودی اختصاص داده می‌شود. بر خلاف صحنه، هر تصویر می‌تواند بیش از یک جسم داشته باشد. تعداد کل اجسام موجود در یک تصویر را با  $K$  و تعداد کل دسته‌های موجود برای اجسام در مجموعه‌داده را با  $O$  نمایش می‌دهیم. فرایند زیر برای هر یک از  $K$  جسم موجود در تصویر اجرا می‌شود:

\* ابتدا یک دسته جسم با توزیع احتمالی  $P(O|E, \pi)$  به تصویر اختصاص داده می‌شود که در آن،  $\pi$  یک ماتریس به ابعاد  $O * E$  و  $\zeta$  یک بردار به طول  $O$  و مقدار اولیه دیریکله پارامتر  $\pi$  است.

\* سپس با داشتن  $O$  می‌توان تمام نواحی  $A$  و  $G$  مرتبط با دسته جسم را تولید نمود. فرایند تولید این نواحی به شکل زیر است:

- متغیر تصادفی مخفی  $z$  که مشخص کننده موضوع است، از توزیع احتمالی  $Mutl(\lambda, |O)$  تولید می‌شود. متغیر  $\lambda$  یک ماتریس به ابعاد  $Z * O$  است که در آن  $Z$  تعداد کل مقادیر مختلف متغیر  $z$  است. به علاوه  $\zeta$  مقدار اولیه دیریکله برای پارامتر  $\lambda$  است.

- نواحی مطلوب از توزیع احتمال  $P(A, G|t, \phi)$  تولید می‌شوند که در آن،  $\phi$  یک ماتریس به ابعاد  $Z * V_o$  است.  $V_o$  تعداد کل کلمات موجود در مجموعه‌داده، به ازای نواحی  $A$  و  $G$  است. پارامتر  $\alpha$  مقدار اولیه دیریکله برای پارامتر  $\phi$  است.

با توجه به متغیرهای تصادفی توضیح داده شده در بالا، توزیع احتمالی توام کل سیستم را می‌توان مطابق با رابطه ۳ تعریف کرد.

$$\begin{aligned} P(E, S, O, X, A, G, t, z, \omega | \rho, \phi, \lambda, \psi, \pi\theta) &= P(E) \cdot P(S|E, \psi) \cdot P(\omega|S, \rho) \\ &\quad \cdot \prod_{m=1}^M P(X_m|t_m, \theta) \cdot P(t_m|\omega) \\ &\quad \cdot \prod_{k=1}^K P(O_k|E, \pi) \\ &\quad \cdot \prod_{n=1}^N P(A_n, G_n|z_n, \phi) \cdot P(z_n|\lambda, O_k) \end{aligned} \tag{۳}$$

به علاوه، با توجه به توضیحات ارائه شده در بالا، هر کدام از عبارات موجود در رابطه ۳ را می‌توان با عبارات معادل آن‌ها که در روابط ۴ تا ۱۰ آمده، جایگزین نمود.

---

<sup>۴۶</sup>Object Image

$$P(S|E, \psi) = Mult(S|E, \psi) \quad (4)$$

$$P(\omega|S, \rho) = Dir(\omega|\rho_j), S = j \quad (5)$$

$$P(t_m|\omega) = Mult(t_m|\omega) \quad (6)$$

$$P(X_m|t_m, \theta) = P(X_m|\theta_j.), t_m = j \quad (7)$$

$$P(O_k|E, \pi) = Mult(O_k|E, \pi) \quad (8)$$

$$P(z_n|\lambda, O_k) = Mult(z_n|\lambda, O_k) \quad (9)$$

$$P(A_n, G_n|z_n, \phi) = P(A_n, G_n|\phi_j.), z_n = j \quad (10)$$

در ک صحنه در این پژوهش، محدود به استخراج سه دسته اطلاعات زیر از تصویر است:

(آ) رخدادی که در تصویر به نمایش گذاشته شده است.

(ب) صحنه‌ای که تصویر در آن ایجاد شده است.

(ج) اجمالی که در تصویر حضور دارند.

با توجه به این محدودیت و با در نظر گرفتن مدل ارائه شده، استفاده از تخمین بیشینه احتمال<sup>۴۷</sup>، می‌تواند برای استخراج اطلاعات مطلوب مفید باشد. از همین رو، تخمین بیشینه احتمال، در سه سطح مختلف (هر سطح برای یک دسته از اطلاعات مطلوب) اعمال می‌شود. در سطح اجسام، احتمال رخداد تصویر ورودی به شرط اجسام موجود مطابق با رابطه ۱۱، احتمال رخداد تصویر ورودی به شرط صحنه، مطابق با رابطه ۱۲ و احتمال رخداد تصویر ورودی به شرط دسته رخداد به نمایش گذاشته شده در تصویر، مطابق با رابطه ۱۳ محاسبه می‌شوند.

$$P(I|O) = \prod_{n=1}^N \sum_j P(A_n, G_n|z_j, O) P(z_j|O) \quad (11)$$

$$P(I|S, \rho, \theta) = \int P(\omega|\rho, S) (\prod_{m=1}^M \sum_{t_m} P(t_m|\omega) P(X_m|t_m, \theta)) d\omega \quad (12)$$

$$P(I|E) \propto \sum_j P(I|O_j) P(O_j|E) P(I|S) P(S|E) \quad (13)$$

فرایند یادگیری این مدل، شامل یافتن بهترین مقادیر برای پارامترهای  $\{\psi, \rho, \pi, \lambda, \theta, \beta\}$  است. این فرایند برای سه پارامتر  $\{\psi, \rho, \theta\}$  با استفاده از روش انتقال پیام متغیر<sup>۴۸</sup> و برای سه پارامتر  $\{\pi, \lambda, \beta\}$  با استفاده از نمونه‌برداری گیبس<sup>۴۹</sup> انجام می‌شود.

آزمایشات انجام شده در این پژوهش، بر روی یک مجموعه‌داده شامل تصاویر از ۸ دسته ورزشی مختلف که در هر دسته، بین ۱۳۷ تا ۲۵۰ تصویر مختلف وجود دارد، انجام شده‌اند. از جمله چالش‌های موجود در این

<sup>۴۷</sup>Maximum Likelihood

<sup>۴۸</sup>Variational Message Passing

<sup>۴۹</sup>Gibbs Sampling

مجموعه‌داده می‌توان به وجود زمینه‌های متعدد و پیچیده در تصاویر، تنوع دسته‌های مختلف اجسام موجود، تنوع اندازه اجسام موجود از یک دسته، تنوع حالت اجسام، تنوع تعداد نمونه‌های یک جسم در یک تصویر و کوچک بودن بیش از اندازه ابعاد اجسام در تصویر اشاره کرد. شکل ۱۰ نمونه‌ای از تصاویر موجود در این مجموعه‌داده را نمایش می‌دهد.

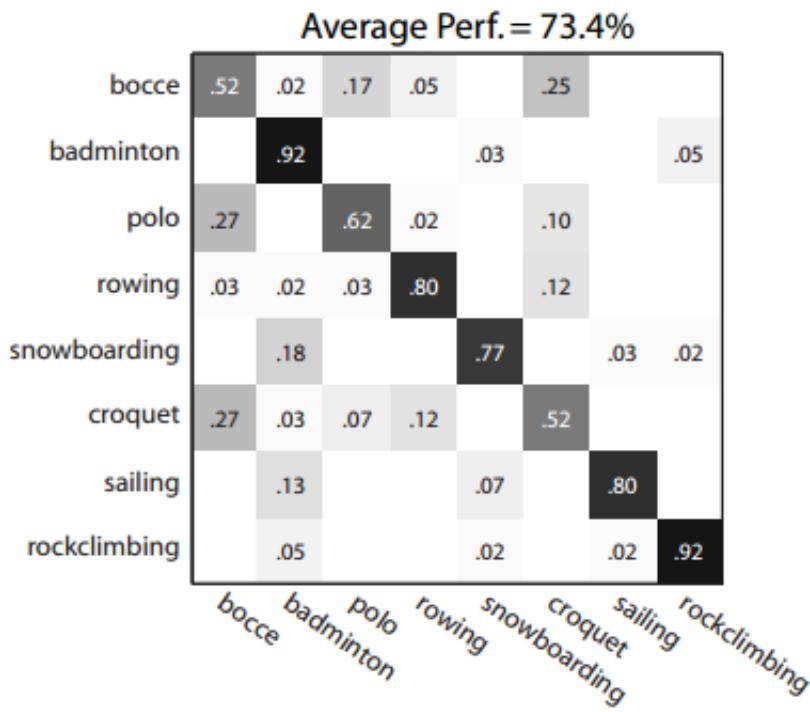


شکل ۱۰: نمونه تصاویر موجود در مجموعه‌داده مورد استفاده [۲]

استفاده از مدل کامل ارائه شده در این پژوهش، منجر به تشخیص صحیح ۷۳.۴٪ از تصاویر شده است. شکل ۱۱ ماتریس درهم‌ریختگی<sup>۵</sup> مربوط به این مدل را نمایش می‌دهد. همان‌طور که در این ماتریس مشخص است، کمترین نرخ تشخیص در بین دسته‌های ورزشی موجود در این مدل، ۵۲٪ و بیشترین نرخ تشخیص

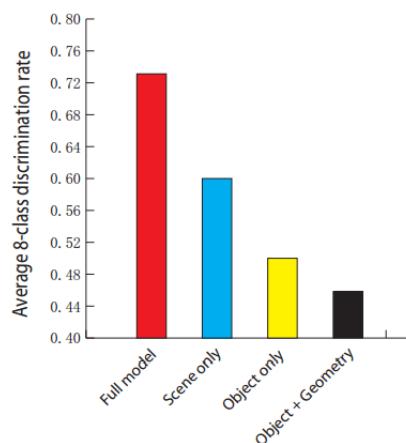
<sup>۵</sup> Confusion Matrix

۹۲٪ است.



شکل ۱۱: ماتریس درهم‌ریختگی مدل کامل ارائه شده برای مجموعه‌داده شامل ۸ دسته تصویر ورزشی. [۲]

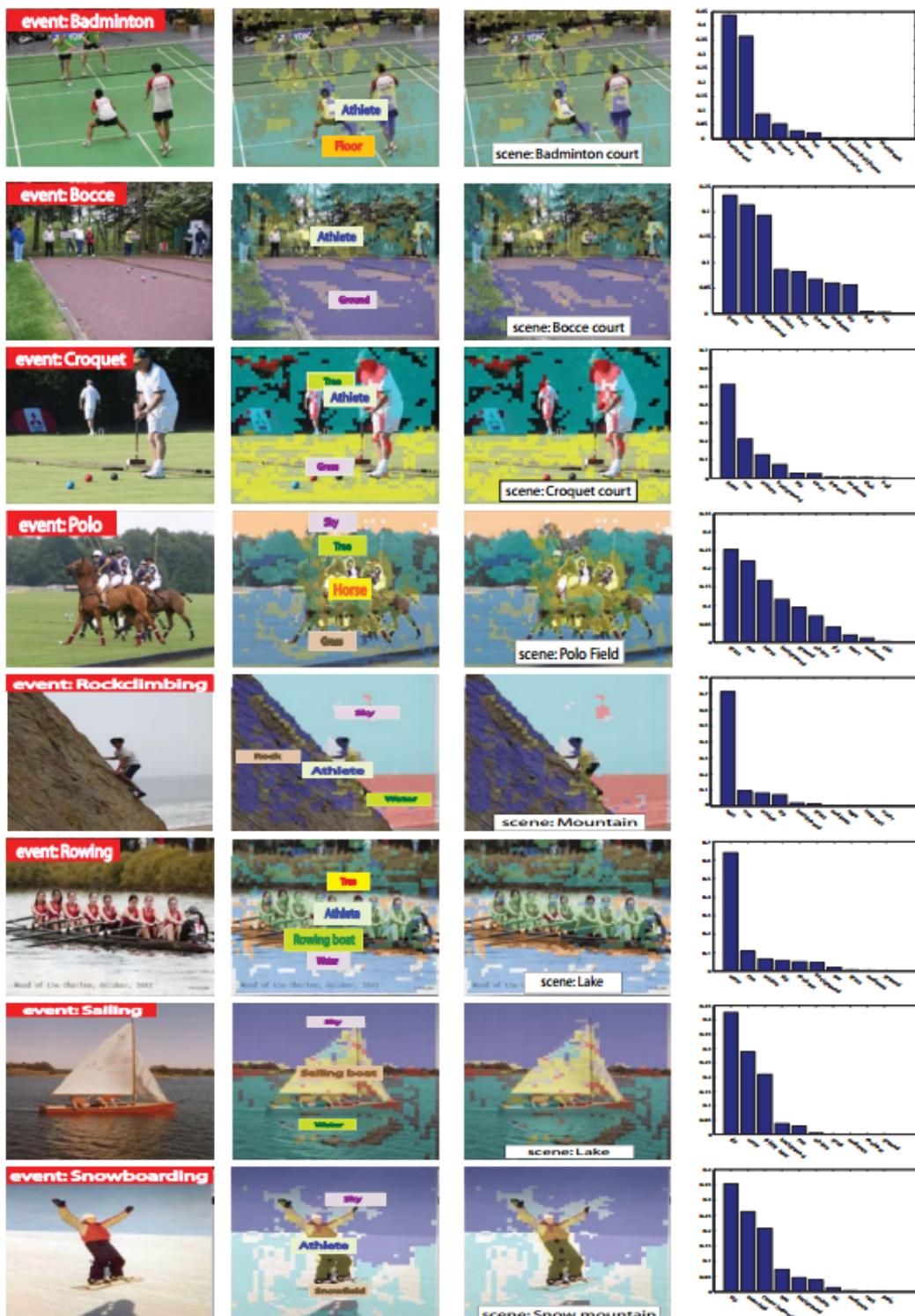
بسته به میزان استفاده از اطلاعات مختلف استخراج شده برای استنتاج، مدل‌های مختلفی به وجود می‌آیند که در شکل ۱۲ نتایج عملکرد هریک از این مدل‌ها با مدل‌های دیگر مقایسه شده است. همان‌طور که در شکل ۱۲ مشخص است، بهترین کارایی مربوط به مدل کامل است. در صورتی که در مدل، فقط از اطلاعات مربوط به صحنه استفاده شود، نتایج بدست آمده اگرچه با نتایج مدل قابل مقایسه نیست، از نتایج مدل مبتنی بر اطلاعات جسم بهتر است.



شکل ۱۲: نتیجه مقایسه مدل‌های مختلف به وجود آمده بسته به سطح اطلاعات مورد استفاده برای استنتاج. [۲]

شکل ۱۳ نتایج نهایی به دست آمده از مدل را نمایش می‌دهد. در این شکل، تصاویر موجود در هر سطر نماینده تصاویر موجود در یکی از دسته‌های ورزشی هستند. ستون اول برچسب به دست آمده از رخداد

موجود در تصویر، ستون دوم برجسب‌های تشخیص داده شده مربوط به اجسام موجود، ستون سوم برجسب اختصاص داده شده مربوط به دسته صحنه و ستون چهارم توزیع مرتب شده اجسام به شرط رخداد را به نمایش می‌گذارند. در نمودارهای موجود در ستون چهارم، محور افقی شامل نام اجسام و محور عمودی مقدار توزیع را نمایش می‌دهد.



شکل ۱۳: نتایج نهایی به دست آمده از مدل بر روی تصاویر. [۲]

## ۴-۲ روش‌های مبتنی بر شبکه‌های عصبی کانولوشنی عمیق

علاوه بر فعالیت‌هایی که در زمینه تولید خودکار شرح بر تصاویر با رویکرد احتمالاتی انجام شده‌اند، تعداد زیادی از پژوهش‌گران تلاش می‌کنند تا با استفاده از روش‌های مبتنی بر شبکه‌های عصبی با این چالش روبرو شوند. در این بخش تعدادی از پژوهش‌هایی را که با استفاده از شبکه‌های عصبی سعی در درک صحنه‌های موجود در تصاویر دارند را مورد بررسی قرار می‌دهیم. شایان ذکر است، در این بخش تنها به بررسی بخشی از پژوهش‌ها که مربوط به درک صحنه است می‌پردازیم و بخش‌هایی از این پژوهش‌ها که مربوط به تولید جملات زبان طبیعی متناسب با تصویر و صحنه درک شده است را در فصل تولید جملات زبان طبیعی بررسی خواهیم نمود.

یکی از مهم‌ترین عملیات‌هایی که به نحوی در تمام پژوهش‌های قبلی وجود داشت، اختصاص یک معنا به قطعه‌های مختلف یک تصویر است. این چالش، در پژوهش‌های مرتبط با تولید خودکار شرح بر تصاویر که با استفاده از روش‌های مبتنی بر شبکه‌های عصبی به دنبال حل مشکل هستند نیز مطرح است. در ابتدا به بررسی یکی از روش‌های اختصاص معنا به هر قطعه از تصویر می‌پردازیم.

### ۱-۴-۲ اختصاص معنا به قطعه‌های مختلف تصویر [۲۰]

در پژوهش [۲۰] روشی ارائه شده است که با استفاده از یک شبکه عصبی کانولوشنی عمیق، علاوه بر این که می‌تواند یک تصویر را به شکل پایین به بالا، در قالب نواحی سلسله‌مراتبی قطعه‌بندی کند، قادر به استفاده به عنوان یک شبکه از پیش آموزش دیده شده در پژوهش‌های مرتبط دیگر باشد.

فرایند تشخیص اجسام در این پژوهش از سه بخش اصلی تشکیل شده است:

۱. طرح پیشنهاداتی برای نواحی به طور مستقل از دسته‌بندی<sup>۵۱</sup>

۲. یک شبکه عصبی عمیق کانولوشنی که وظیفه استخراج ویژگی برای هر ناحیه را بر عهده دارد (طول بردار ویژگی استخراج شده برای تمام نواحی یکسان است).

۳. مجموعه‌ای از ماشین‌های بردار پشتیبان خطی مخصوص هر دسته

در ادامه به بررسی نحوه پیشنهاد نواحی و شبکه عصبی کانولوشنی عمیق مورد استفاده در ای پژوهش می‌پردازیم.

#### ۱. طرح پیشنهاد نواحی

روش‌های مختلفی برای پیشنهاد نواحی ارائه شده‌اند که در اینجا از روشی موسوم به جستجوی انتخابی<sup>۵۲</sup> استفاده می‌شود. نسخه‌های مختلفی از این روش ارائه شده است. نسخه ارائه شده در پژوهش [۲۱]، یکی از سریع‌ترین نسخه‌های ارائه شده است که در این بخش از همین روش استفاده می‌شود.

در پژوهش [۲۱] دو ویژگی مطرح شده است که یک جستجوی انتخابی برای ارائه نواحی معنایی تصویر باید آن‌ها را داشته باشد. ویژگی اول این است که اجسام موجود در فضای می‌توانند در هر اندازه‌ای باشند و در نتیجه نواحی ارائه شده باید بتوانند ابعاد مختلف داشته باشند. این ویژگی عموماً با روش‌های سلسله‌مراتبی

<sup>۵۱</sup>Category-independent region proposals

<sup>۵۲</sup>Selective Search

قابل دستیابی است. ویژگی دوم این است که نواحی مختلف باید براساس ویژگی‌های مختلفی تولید شوند. در صورتی که یک ویژگی مثل رنگ، بافت، روشنایی یا مواردی از این دست، به عنوان تنها ویژگی برای تشخیص نواحی به کار گرفته شود، الگوریتم قادر به ارائه نواحی مناسب در شرایط مختلف نخواهد بود. بنابراین ترکیب چند معیار و ویژگی باید برای تشخیص نواحی مورد استفاده قرار بگیرد.

برای دستیابی به ویژگی اول، ابتدا نواحی اولیه کوچکی روی تصویر ایجاد می‌شود. سپس با اتخاذ یک روش حریصانه و تعریف یک معیار شباخت بین نواحی همسایه، ناحیه‌هایی که شباخت زیادی با یکدیگر دارند و همسایه هستند، با هم ترکیب شده و یک ناحیه بزرگ‌تر ساخته می‌شود. به این ترتیب یک روش سلسله‌مراتبی برای ساخت نواحی با ابعاد مختلف به دست می‌آید. برای دستیابی به ویژگی دوم، از فضاهای رنگی مختلف، معیارهای شباخت مختلف و نواحی اولیه متفاوت و ترکیب پاسخ این ویژگی‌ها با هم برای ارائه نواحی و ترکیب نواحی کوچک‌تر استفاده می‌شود.

## ۲. شبکه عصبی کانولوشنی عمیق (استخراج ویژگی‌ها)

در این بخش از یک شبکه عصبی کانولوشنی عمیق از پیش‌آموزش دیده برای استخراج ویژگی از هر ناحیه ارائه شده در قسمت قبل، استفاده می‌شود. بردار ویژگی استخراج شده برای هر ناحیه یک بردار شامل ۴۰۹۶ مولفه است که خروجی شبکه کریشفسکی<sup>۵۳</sup> آزمایش شده در چالش دسته‌بندی اجسام مسابقه ImageNet است. اطلاعات دقیق درباره این شبکه عصبی در پژوهش [۲۲] در دسترس است.

شبکه عصبی کانولوشنی عمیق ارائه شده در این پژوهش با استفاده از یک مجموعه‌داده<sup>۵۴</sup> آموزش دیده شده است. از این شبکه عصبی که تحت عنوان RCNN<sup>۵۵</sup> شناخته می‌شود می‌توان به عنوان یک شبکه از پیش‌آموزش دیده استفاده کرد.

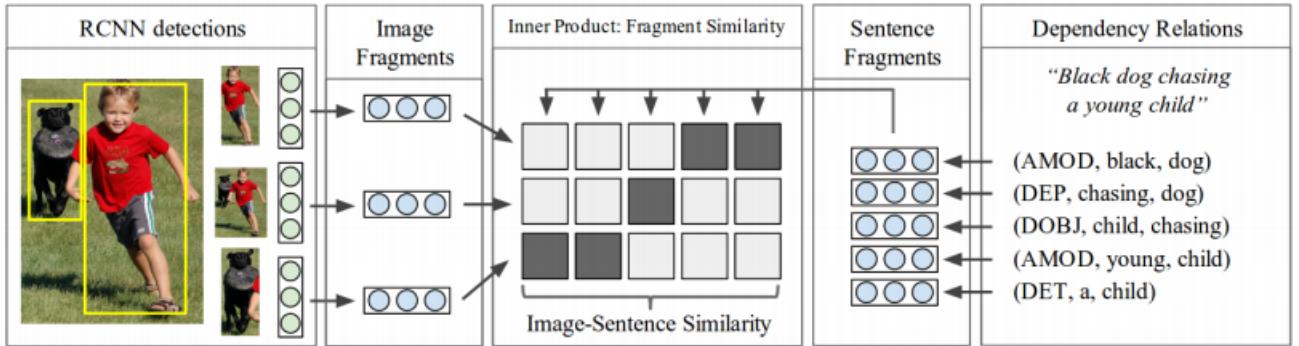
### ۲-۴-۲ ناحیه‌بندی عمیق تصاویر به منظور نگاشت دوطرفه جملات و تصاویر [۲۳]

مدل ارائه شده در این پژوهش، مدلی است که قادر به نگاشت دوطرفه تصاویر و جملات به یکدیگر است. شکل ۱۴ طرح‌واره‌ای از این مدل را نمایش می‌دهد. ورودی مدل در سمت چپ، تصاویر و در سمت راست، جملات هستند. در این مدل، ابتدا تصاویر ورودی با استفاده از یک شبکه عصبی RCNN تبدیل به نواحی مختلف شده و برای هر ناحیه یک بردار ویژگی ۴۰۹۶ بعدی استخراج می‌شود. سپس با اعمال روش خاصی روی جملات ورودی از سمت راست (که در بخش تولید جملات زبان طبیعی به بررسی آن خواهیم پرداخت) قطعات مختلف موجود در جملات نیز استخراج شده و بین هر قطعه از جمله با تمام نواحی استخراج شده از تصویر یک معیار شباخت محاسبه می‌شود و شبیه‌ترین قطعه جمله با ناحیه مربوط به خود در تصویر، جفت می‌شوند.

<sup>۵۳</sup>Krizhevsky

<sup>۵۴</sup>ILSVRC 2012

<sup>۵۵</sup>Regional Convolutional Neural Network



شکل ۱۴: مدل استفاده شده برای نگاشت دوطرفه تصاویر و جملات به یکدیگر با استفاده از شبکه عصبی عمیق کانولوشنی. [۲۳]

در این پژوهش پس از ناحیه‌بندی تصویر توسط شبکه RCNN، برای هر تصویر ۱۹ ناحیه استخراج می‌شود. این ۱۹ ناحیه در کنار تصویر اصلی، یک مجموعه شامل ۲۰ تصویر ایجاد می‌کنند که در پردازش‌های بعدی مورد استفاده قرار خواهند گرفت. در این مرحله باید تمام تصاویر موجود را با استفاده از یک نگاشت به فضای برداری ویژگی‌ها تبدیل نمود. برای این کار از رابطه  $I_b$  استفاده می‌شود. در این رابطه،  $RCNN_{\theta_c}(I_b)$  مجموعه تمام پیکسل‌های موجود در ناحیه  $b$ ، شبکه عصبی آموزش‌دیده است که در آن  $\theta_c$  مجموعه پارامترهای بهینه موجود در شبکه است. بردار حاصل  $\nu_i$  برای تصویر  $i$ ، بردار نگاشت تصویر به فضای معنایی خواهد بود که محاسبه مقادیر آن مبتنی بر پیشنهاد نواحی معنایی مختلف و محاسبه ویژگی‌های مختلف روی هر ناحیه است.

$$\nu = W_m[RCNN_{\theta_c}(I_b)] + b_m \quad (14)$$

از طرفی با در نظر گرفتن بردار  $s_j$  به عنوان بردار حاصل از نگاشت جمله  $z$  به فضای معنایی و در نظر گرفتن ضرب داخلی به عنوان شباهت،  $s_j^T \cdot \nu_i$  معیار شباهت بین یک تصویر و یک جمله را تعریف می‌کند. با توجه به توضیحات ارائه شده، می‌توان تابع هدف را برای شبکه کلی معادل سیستم ارائه داد. دو هدف اصلی در این شبکه قابل تعریف است:

۱. رتبه‌بندی سراسری تصاویر و جملاتی که در فرایند محاسبات شبکه عصبی بیشترین شباهت را با یکدیگر دارند باید در واقعیت هم بیشترین شباهت و ارتباط را داشته باشند.
۲. هم ترازسازی ناحیه‌ای<sup>۵۶</sup> نواحی استخراج شده تصویر و عبارات استخراج شده جملات که در محاسبات شبکه عصبی بیشترین شباهت را با یکدیگر دارند، باید در واقعیت هم بیشترین شباهت و ارتباط را داشته باشند.

با توجه به مطالب گفته شده، می‌توان تابع هدف کلی را مطابق با رابطه ۱۵ تعریف کرد. در این رابطه،  $\Theta$  مجموعه پارامترهای شبکه عصبی شامل  $\{W_m, b_m, \theta_c, W_e, W_R\}$  است (پارامترهای  $W_e$  و  $W_R$  مربوط به بخش

<sup>۵۶</sup>Fragment Alignment

تحلیل جمله هستند که در فصل مربوطه بررسی خواهند شد).  $C_F$  تابع هدف همترازسازی ناحیه‌ای،  $C_G$  تابع هدف سراسری،  $\alpha$  و  $\beta$  دو ابرپارامتر<sup>۵۷</sup> (با آزمون و خطا تعیین می‌شوند) و  $\|\cdot\|_2^2$  یک عبارت تنظیم‌کننده<sup>۵۸</sup> هستند.

$$C(\Theta) = C_F(\Theta) + \beta C_G(\Theta) + \alpha \|\Theta\|_2^2 \quad (15)$$

در ادامه به تعریف هریک از اهداف بیان شده می‌پردازیم.

## ۱. همترازسازی ناحیه‌ای

هدف از همترازسازی ناحیه‌ای این است که اگر عبارتی از یک جمله با یک تصویر شباهت زیادی پیدا کرد، حداقل یک ناحیه از تصویر وجود داشته باشد که نمایش‌دهنده این عبارت باشد و بقیه نواحی تصویر، ارتباط کمی با این عبارت داشته باشند. به عبارت بهتر، در صورتی که شباهت یک عبارت از یک جمله با یک تصویر از حدی بیشتر شد، شباهت حداقل یکی از نواحی موجود در تصویر با این عبارت زیاد شده و شباهت بقیه نواحی تصویر با آن کم شود. این فرض در سه حالت، رد می‌شود. اولین حالت، حالتی است که در آن ناحیه‌ای که در واقعه نمایش‌دهنده عبارت است، توسط RCNN تشخیص داده نشده باشد. دومین حالت، حالتی است که عبارت موجود به هیچ بخشی از ویژگی‌های بصری تصویر اشاره نکند و آخرین حالت، حالتی است که عبارت توصیف‌کننده، در هیچ یک از تصاویر دیگر تکرار نشده باشد در صورتی که ممکن است تصاویر دیگری هم وجود داشته باشند که شامل ویژگی‌های بصری متناظر با عبارت باشند. با توجه به شرایطی که فرض در آن‌ها نقض می‌شود، می‌توان آن را یک فرض خوب تلقی کرد که در اکثر موارد عملکرد خوبی دارد.

رابطه ۱۶ تابع هدف همترازسازی ناحیه‌ای را تعریف می‌کند. در این رابطه،  $y_{ij}$  برای تصویر  $i$ ام و جمله  $j$ ام در صورتی که با هم در مجموعه‌داده حضور داشته باشند،  $+1$  و در غیر این صورت،  $-1$  خواهد شد.

$$C_{\circ}(\Theta) = \sum_i \sum_j \max(0, 1 - y_{ij} \nu_i^T \cdot s_j) \quad (16)$$

تابع  $C_{\circ}$  تعریف شده، باعث می‌شود در حالاتی که تصویر و عبارت، در مجموعه‌داده، با یکدیگر وارد شده باشند امتیاز تابع هدف بیشتر از  $+1$  شود و در غیر این صورت از  $-1$  کمتر شود. شکل ۱۵، دو نمونه از تصاویر و جملات موجود در مجموعه‌داده را نمایش می‌دهد.  $C_{\circ}$  در سلول‌هایی که با رنگ قرمز مشخص شده‌اند، امتیاز را به سمت کمتر از  $-1$  حرکت می‌دهد و در بقیه سلول‌ها به سمت بیشتر از  $+1$ .

به عبارت بهتر،  $C_{\circ}$  یک امتیاز برای مجموع تفاوت‌های نواحی مختلف از تصاویر با عبارات مختلف جملات است. به دلیل این‌که این معیار، باعث دیده نشدن موارد کمیاب می‌شود، با متغیر گرفتن پارامتر  $y_{ij}$  سعی در یافتن کمترین مقدار آن می‌کنیم. رابطه ۱۷ معیار متناظر با هدف کلی همترازسازی ناحیه‌ای را بیان می‌کند.

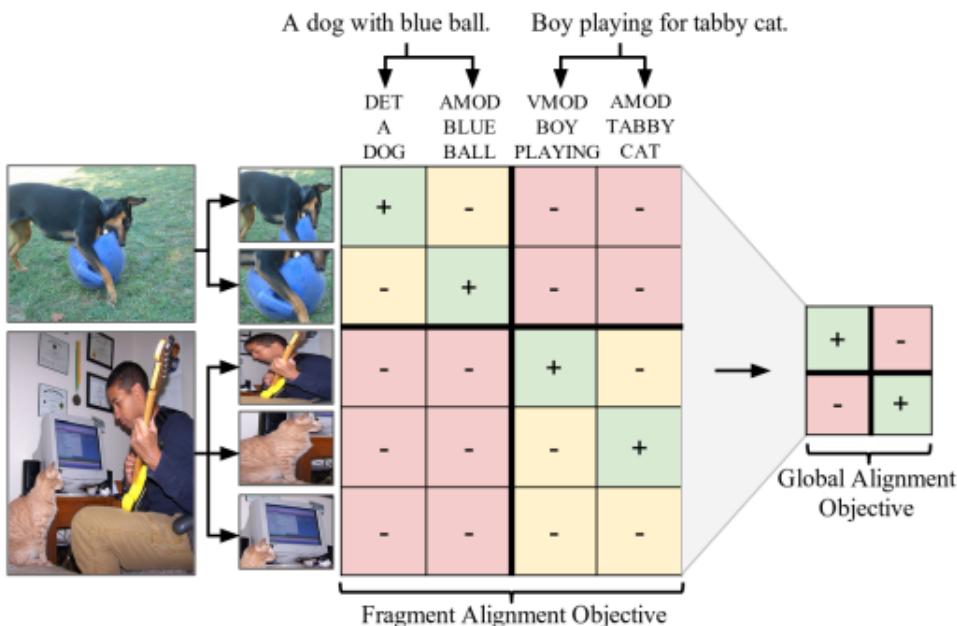
<sup>۵۷</sup>Hyperparameter

<sup>۵۸</sup>Regularization Term

$$C_F(\Theta) = \min_{y_{ij}} C_\circ(\Theta)$$

$$\text{s.t. } \sum_{i \in p_j} \frac{y_{ij} + 1}{2} \geq 1 \wedge \quad y_{ij} = -1, \forall i, j; m_\nu(i) \neq m_s(j) \wedge y_{ij} \in \{+1, -1\} \quad (17)$$

در این رابطه،  $p_j$  مجموعه تصاویر موجود در کیسه مثبت<sup>۵۹</sup> مربوط به عبارت زام است. شایان ذکر است، تنها تصاویری که در مجموعه‌داده همراه با عبارت زام مشاهده شده‌اند در کیسه مثبت مربوط به این عبارت قرار می‌گیرند و بقیه تصاویر در کیسه منفی<sup>۶۰</sup> این عبارت قرار می‌گیرند.  $m_\nu(i)$  و  $m_s(j)$  به ترتیب، شماره تصویر و عبارت را در مجموعه‌داده مشخص می‌کنند.



شکل ۱۵: دو نمونه از تصاویر و جملات مرتبط با آن‌ها و نتایج عملکرد اهداف تعریف شده روی آن‌ها. سطرها نمایش‌دهنده نواحی مختلف تصویر و ستون‌ها نمایش دهنده قطعه‌های مختلف جملات هستند. سلول‌های قرمز رنگ حالتی هستند که در آن‌ها  $y_{ij} = 1$ ، سلول‌های زرد نمایش‌دهنده اعضای کیسه‌های مثبت هستند که در آن‌ها  $y_{ij} = -1$  است. [۲۳]

## ۲. رتبه‌بندی سراسری

هدف از رتبه‌بندی سراسری این است که شباهت بین یک تصویر و یک جمله، بیشینه شود اگر و تنها اگر تصویر و جمله در واقعیت نیز بیشترین شباهت را به یکدیگر داشته باشند. برای این منظور، ابتدا یک امتیاز شباهت بین یک تصویر و یک جمله تعریف می‌شود. این امتیاز مطابق با رابطه ۱۸ تعریف شده و برابر است با میانگین امتیاز شباهت دوبه‌دوی نواحی مختلف تصویر با عبارات مختلف جمله.

$$S_{kl} = \frac{1}{|g_k|(|g_l| + n)} \sum_{i \in g_k} \sum_{j \in g_l} \max(\circ, \nu_i^T \cdot s_j) \quad (18)$$

<sup>۵۹</sup>Positive Bag

<sup>۶۰</sup>Negative Bag

از آن جا که برای دسته‌بندی از روش mi\_SVM استفاده می‌شود، تمام امتیازها به صفر محدود می‌شوند. مقدار  $n$  که در مخرج کسر اضافه شده است، به صورت تجربی و با آزمون و خطا به دست آمده که نتایج را بهبود می‌بخشد. مقدار پیشنهاد شده در پژوهش،  $n = 5$  است. تابع کلی هدف سراسری مطابق با رابطه ۱۹ تعریف می‌شود.

$$C_G(\Theta) = \sum_k (\sum_l \max(0, S_{kl} - Skk + \Delta) + \sum_l \max(0, S_{lk} - Skk + \Delta)) \quad (19)$$

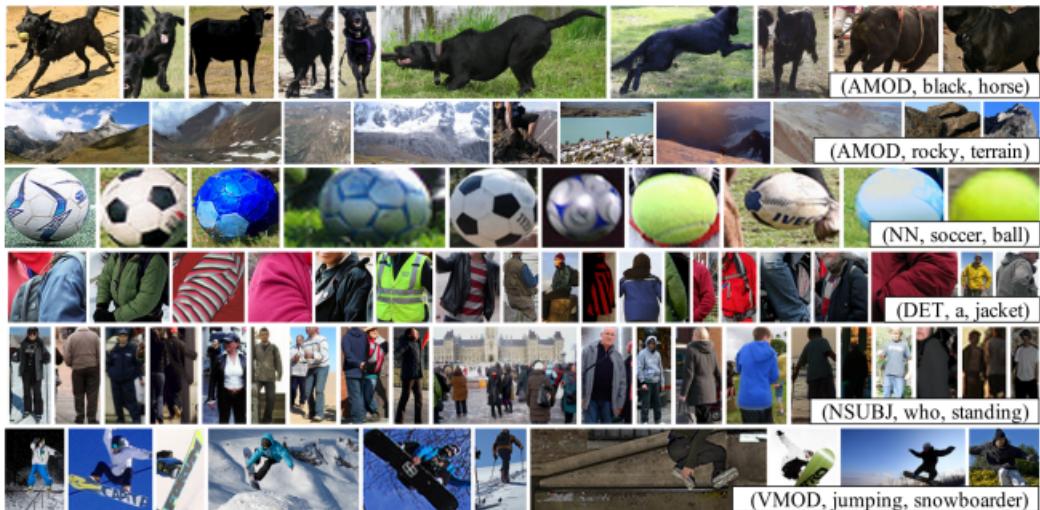
در رابطه ارائه شده،  $\Delta$  یک ابرپارامتر است که با آزمون و خطا به دست می‌آید. عبارت اول درون پرانتز بیان کننده امتیاز تصویر و عبارت دوم بیان کننده امتیاز جمله هستند.

شکل ۱۶ نتایج روش پیشنهاد شده در این پژوهش را ارائه می‌دهد. همان‌طور که در شکل مشخص است، این شبکه قادر به تشخیص اجسام مختلف در تصویر و تولید یک سه‌تایی متناظر هر جسم (ناحیه معنایی) مبتنی بر جملات موجود در مجموعه‌داده مورد استفاده است.



شکل ۱۶: نتایج نهایی شبکه عصبی ارائه شده. برای هر ناحیه معنایی از تصویر، یک سه‌تایی مبتنی بر جملات موجود در مجموعه‌داده تولید شده است. همین‌طور ۵ جمله تولید شده برای هر تصویر به ترتیب امتیاز، درج شده‌اند. [۲۳]

به علاوه، با توجه به مدل ارائه شده و نگاشت دوطرفه موجود بین تصاویر و جملات، می‌توان با ورودی دادن یک جمله، تصاویر مربوط به آن جمله را استخراج نمود. شکل ۱۷ با ثابت در نظر گرفتن جملات، تصاویر مربوط به هر جمله را استخراج و نمایش داده است. هر سطر از این شکل، نمایش‌دهنده تصاویر استخراج شده مرتبط با جمله موجود در آن سطر است.



شکل ۱۷: نتایج حاصل از جستجوی جملات، با ورودی دادن یک جمله، شبکه عصبی ارائه شده در این پژوهش، قادر به استخراج تصاویر مربوط به آن جمله است. [۲۳]

روش ارائه شده در این پژوهش، به طور کامل و دقیق در پژوهش [۲۴] هم مورد استفاده قرار گرفته است، با این تفاوت که در فرایند تحلیل جمله، تغییراتی ایجاد شده است. جزئیات این روش در فصل تولید جملات زبان طبیعی مورد بررسی قرار خواهد گرفت.

## ۲-۵ جمع‌بندی

اولین مرحله از فرایند تولید خودکار شرح برای تصاویر، مرحله درک صحنه است. در این مرحله، تصاویر ورودی تحت عملیات مختلفی به فضای معنایی نگاشت می‌شوند. فضای معنایی در اینجا، می‌تواند فضای شامل میدان‌های اطلاعاتی از پیش تعیین شده (مانند فضای سه‌تایی‌های «جسم، رخداد، صحنه») یا فضای بردار ویژگی‌ها باشد. روش‌های مختلفی برای نگاشت تصویر ورودی به فضای معنایی ارائه شده است که به طور کلی می‌توان عموم آن‌ها را به دو بخش تقسیم کرد:

### ۱. روش‌های مبتنی بر مدل‌های گرافی احتمالاتی

در این روش‌ها با استفاده از مدل‌های استاندارد گرافی احتمالاتی موجود یا با ارائه یک مدل گرافی احتمالاتی، تصویر ورودی به فضای معنایی نگاشت می‌شود. در روش‌های مبتنی بر این مدل‌ها، با ارائه یک توزیع احتمال برای نقاط مختلف در فضای معنایی، محتمل‌ترین نقطه برای تصویر به عنوان نقطه نظر تصویر، انتخاب می‌شود.

#### (آ) مدل میدان تصادفی مارکف

یک نمونه از روش‌های مبتنی بر مدل میدان تصادفی مارکف که برای درک صحنه از آن استفاده شده است، در پژوهش [۱۳] ارائه شده است. درک صحنه در این پژوهش با ارائه یک سه‌تایی «جسم،

فعالیت، صحنه» به ازای هر تصویر، تعریف شده است. مبتنی بر همین تعریف، یک مدل میدان تصادفی مارکف شامل سه گره که دوبهدو به هم متصل هستند، تعریف شده است. هر یک از گرههای موجود در این مدل، نماینده یکی از میدان‌های سه‌گانه تعریف شده در فضای معنایی هستند. با تعریف توابع پتانسیل مختلف روی هر گره و توابع پتانسیل مختلف روی هر یال، یک تابع توزیع توام برای تمام متغیرهای تصادفی موجود در مدل ارائه شده است.

با محاسبه مقادیر پتانسیل برای تصاویر مختلف موجود در مجموعه آموزشی و با استفاده از یک ماشین بردار پشتیبان، بردارهای ویژگی شاخص برای هر گره محاسبه می‌شوند. از این بردارهای ویژگی بعده برای انطباق تصاویر با مقادیر مختلف در هر گره استفاده می‌شود.

در این پژوهش، با یافتن نزدیک‌ترین همسایه‌های یک تصویر بر حسب معیار شباهت با بردارهای ویژگی شاخص و میانگین‌گیری روی مقادیر هر گره، بهترین انطباق تصویر و نقاط فضای معنایی به دست می‌آید. به این ترتیب، برای هر تصویر ورودی، می‌توان نقطه نظری در فضای معنایی را مشخص کرد.

#### (ب) مدل میدان تصادفی شرطی

در پژوهش [۱۲] یک مدل میدان تصادفی شرطی سلسله‌مراتبی برای درک صحنه ارائه شده است که شامل دو سطح انتزاع است. برای گرههای موجود در هریک از سطوح انتزاع مدل، یک دسته متغیر تصادفی تعریف شده و برای کل مدل سه نوع تابع پتانسیل مختلف معرفی شده است.

اولین دسته از توابع پتانسیل معرفی شده در این بخش، توابع پتانسیل قطعه‌بندی یگانی هستند که به منظور یکپارچه‌سازی نقاط داخل یک قطعه تعریف شده‌اند. توابع پتانسیل دیگری برای انطباق بین متغیرهای تصادفی موجود در بین دو سطح انتزاع تعریف شده‌اند که در صورت مغایرت مقادیر اختصاص داده شده به متغیرهای موجود بین دو سطح، مقدار  $\lambda$  – و در غیر این صورت مقدار صفر دارند. این توابع در شرایطی که مقادیر متغیرهای موجود در دو سطح با هم یکسان نباشد، یک مقدار جریمه به تابع هدف اضافه می‌کنند. آخرین دسته از توابع پتانسیل مورد استفاده، برای انطباق تصویر با دسته تشخیص داده شده اجسام تعریف شده است که توسط فلزنسوالب ارائه شده و به روش دی‌پی‌ام مشهور است.

(ج) سایر مدل‌های گرافی احتمالی در پژوهش [۲]، یک مدل گرافی احتمالی مولد برای نگاشت تصویر به فضای معنایی ارائه شده است. در این مدل، از دو سطح تصویر استفاده شده است؛ تصویر سطح جسم و تصویر سطح صحنه. برای تصویر سطح صحنه، یک متغیر تصادفی، بیان‌کننده دسته صحنه و برای تصویر سطح جسم دو متغیر تصادفی، بیان‌کننده دسته و شکل جسم، ارائه شده است. روابط بین متغیرهای تصادفی در این پژوهش، براساس نحوه تولید متغیرهای تصادفی و روابط منطقی موجود بین آن‌ها طراحی شده‌اند.

تصویر ورودی در این پژوهش، ابتدا به نواحی کوچک  $10 \times 10$  تقسیم می‌شود و مطابق با روش توضیح داده شده، مقدار توابع پتانسیل مختلف برای هر کدام از متغیرهای تصادفی، در هر ناحیه، محاسبه می‌شود. در این پژوهش، یک تابع احتمال شرطی برای متغیرهای تصادفی ارائه شده است که در مرحله

استنتاج، با استفاده از روش تخمین بیشترین احتمال، برچسب‌های هر تصویر مشخص می‌شوند.

## ۲. روش‌های مبتنی بر استفاده از شبکه‌های عصبی کانولوشنی عمیق

در این روش‌ها، با ارائه یک شبکه عصبی کانولوشنی عمیق و تعریف کردن تابع هدف برای شبکه، تابع نگاشت تصویر و فضای معنا تشکیل می‌شود. پس از ارائه تابع هدف برای هر شبکه، با بهینه‌سازی آن تابع، پارامترهای موجود در شبکه آموزش داده می‌شوند.

در پژوهش [۲۰]، روشی ارائه شده است که طی آن یک تصویر، به نواحی کوچک‌تر تقسیم می‌شود به طوری که هر ناحیه به وجودآمده، به طور یکپارچه، حاوی یک جسم باشد و هر جسم تنها در یک ناحیه قرار بگیرد. این روش موسوم به روش RCNN است. در این روش، دو ویژگی برای یک ناحیه‌بندی خوب در تصاویر ارائه شده است و پیرو این ویژگی‌ها، روشی برای طرح نواحی پیشنهادی در یک تصویر که دارای این دو ویژگی باشد، ارائه شده است.

ویژگی مطرح شده اول برای ناحیه‌بندی تصاویر این است که، ناحیه‌های ایجاد شده در هر تصویر، می‌توانند در ابعاد مختلف وجود داشته باشند زیرا اجسام موجود در تصاویر، ممکن است اندازه و تعداد متفاوتی داشته باشند. دومین ویژگی برای یک ناحیه‌بندی خوب، این است که معیار انتخاب نواحی نباید برای تمام تصاویر، یکسان در نظر گرفته شود؛ زیرا معیارهای مختلف برای ناحیه‌بندی تصاویر در شرایط مختلف، رفتارهای متفاوتی از خود نشان می‌دهند. بنابراین باید از معیارهای مختلف برای تعیین نواحی استفاده نمود.

در این پژوهش، ابتدا تصاویر مطابق با یک معیار اولیه، به مجموعه‌ای از نواحی اولیه تقسیم می‌شوند. سپس با استفاده از معیارهای مختلف مانند فضاهای رنگی مختلف، معیارهای شباht مختلف و نقاط اولیه متفاوت، با پیروی از یک روش حریصانه، نواحی کوچک‌تر که به یکدیگر شبیه‌تر هستند با هم ترکیب شده و نواحی بزرگ‌تر را می‌سازند. نواحی ایجاد شده در این روش، سپس به یک شبکه عصبی کانولوشنی عمیق داده می‌شوند و برای هر ناحیه، یک بردار ویژگی ۴۰۹۶ بعدی ایجاد می‌شود که هر ناحیه با آن بازنمایی شود.

در پژوهش [۲۳] با استفاده از روش RCNN و تعریف دو تابع هدف دیگر برای شبکه، روشی ارائه شده است که طی آن بتوان تصاویر و جملات را به طور دوطرفه به یکدیگر نگاشت کرد. تابع هدف تعریف شده در این پژوهش، دو تابع مختلف هستند. اولین تابع هدف، یک تابع هدف سراسری است. این تابع به این منظور تعریف شده است که تصاویر و جملاتی که مطابق با محاسبات شبکه عصبی ارائه شده، بیشترین شباهت را با یکدیگر دارند، در واقعیت هم شبیه‌ترین تصاویر و جملات به یکدیگر باشند. تابع هدف دوم برای این شبکه به این شکل تعریف شده است که نواحی استخراج شده از تصویر و عبارات استخراج شده از جملات که در روش ارائه شده، بیشترین شباهت را به یکدیگر دارند، در واقعیت هم بیشترین شباهت و ارتباط را با یکدیگر داشته باشند.

در این پژوهش، تصاویر ورودی با استفاده از روش RCNN به نواحی مختلف تقسیم شده و ۱۹ ناحیه با بیشترین اطمینان از بین این نواحی انتخاب می‌شود. این ۱۹ ناحیه به همراه خود تصویر به عنوان ۲۰ تصویر مختلف مورد استفاده قرار می‌گیرند. جملات ورودی با استفاده از روشی که در فصل تولید جملات زبان طبیعی توضیح داده خواهد شد، به عبارات مختلف تقسیم می‌شوند و بین هر عبارت استخراج شده و هر یک از ۲۰ تصویر موجود، یک معیار شباهت محاسبه شده و بیشترین شباهت‌ها با هم درنظر گرفته می‌شوند.

معیار شباهت مورد استفاده در این روش، ضرب داخلی بین بردارهای ویژگی عبارات و نواحی است. عبارات و نواحی که بیشترین شباهت را با یکدیگر دارند برای تولید جمله به مرحله بعد، ارسال می‌شوند.

## ٣ فصل سوم

### تولید جمله

## ۱-۳ تولید جمله

در این بخش، به بررسی چالش تولید جمله و روش‌های پیشنهاد شده برای حل این چالش، در طول زمان، خواهیم پرداخت. در ابتداء ایده‌های اولیه در این مسیر را بیان نموده و به طور خلاصه مورد بحث قرار خواهیم داد و در ادامه به بررسی تفصیلی روش‌های مبتنی بر استفاده از شبکه‌های عصبی بازگشتی در تولید جملات زبان طبیعی، می‌پردازیم.

چالش تولید جمله، متوجه ساخت جملاتی به زبان طبیعی است، به طوری که از لحاظ دستور زبان، املا و معنا صحیح باشند. از طرفی با توجه به هدف اصلی ما که تولید شرح بر تصاویر است، جملات تولید شده باید علاوه بر این که شرط صحت مذکور را ارضا می‌کنند، با تصویر ورودی، صحنه توصیف شده در تصویر و رخداد به نمایش کشیده شده، هم‌خوانی داشته باشند. تضمین این هم‌خوانی از جمله معضلات دیگری است که باید برای آن چاره‌ای اندیشید.

مساله تولید خودکار جملات زبان طبیعی، یکی از مسائلی است که از دیرباز در هوش مصنوعی مطرح بوده و دارای کاربردهای فراوانی است. به عنوان یک تعریف دقیق از این مساله، می‌توان به «فرایند تولید جملات زبان طبیعی با استفاده از داده‌های غیر قابل تفسیر برای کاربران عادی در جهت افزایش قابلیت تفهیم داده به آن‌ها [۲۵]» اشاره کرد. داده‌های اولیه که جملات با استفاده از آن‌ها تولید می‌شوند، می‌توانند شامل انواع داده‌های غیر متنی از جمله نمودارها، تصاویر، اعداد و مواردی از این دست باشند. در ادامه، ابتدا کاربردهای مساله تولید خودکار جملات زبان طبیعی<sup>۶۱</sup>، خواهیم پرداخت.

## ۲-۳ کاربردها [۲۵]

کاربردهای بسیار زیاد و متنوعی برای تولید خودکار جملات زبان طبیعی با استفاده از داده‌های غیر متنی، ارائه شده است. در این بخش، برای روشن شدن اهمیت این مساله، تعدادی از این کاربردها را بیان خواهیم نمود.

۱. تولید خودکار شرح بر پیش‌بینی وضع آب و هوا با استفاده از نقشه‌های گرافیکی آب و هوا
۲. تولید خلاصه‌ای درباره داده‌های آماری استخراج شده از یک پایگاه داده
۳. توصیف یک زنجیره استدلالی، منتج از فرایند تصمیم‌گیری یک سیستم خبره
۴. تولید پاسخ برای پرسش‌ها در مورد یک جسم در یک سامانه مبتنی بر دانش

<sup>۶۱</sup>Natural Language Generation (NLG)

موارد ذکر شده در بالا، تنها نمونه‌ای از کاربردهای وسیع این مساله در زندگی‌های روزمره را نمایش می‌دهد.

### ۳-۳ روش تولید زبان طبیعی

یکی از حوزه‌های پویا در زمینه هوش مصنوعی و پردازش متن، حوزه تولید زبان طبیعی است. این حوزه شامل فعالیت‌هایی است که به تولید جمله زبان طبیعی متناظر با داده‌های قابل تفسیر برای ماشین مانند پایگاه‌های دانش<sup>۶۲</sup> و یا قالب‌های منطقی<sup>۶۳</sup> و مواردی از این دست، می‌پردازد. روش‌های کلی که در چارچوب کاری<sup>۶۴</sup> پژوهش‌های این زمینه مورد استفاده قرار می‌گیرد، به طور کلی شامل مراحل زیر است [۲۵].

#### ۱. برنامه‌ریزی متن<sup>۶۵</sup>

در این مرحله، ابتدا محتوا مورد نیاز، انتخاب می‌شود و چارچوب کلی برای کل متن، طرح ریزی می‌شود. انتخاب محتوا و طرح ریزی چارچوب کلی متن، با تکیه بر کاربرد مورد نظر در پژوهش و داده‌هایی که نیاز به تفسیر زبانی دارند، انجام می‌شود.

#### ۲. برنامه‌ریزی جمله<sup>۶۶</sup>

در این مرحله، کلمات مورد نیاز انتخاب می‌شوند، عبارات زبانی مناسب تولید می‌شوند و به طور دقیق کنارهم قرار می‌گیرند تا جملات را تشکیل دهند. انتخاب کلمات و ساخت عبارات در این بخش، بر اساس طرح کلی پی‌ریزی شده برای متن در مرحله قبل، انجام می‌شود. کنارهم قرار دادن عبارات زبانی مورد نیاز نیز، با استفاده از قواعد دستور زبان که عموماً به شکل پایگاه دانش موجود است، انجام می‌گردد.

#### ۳. تحقق زبانی<sup>۶۷</sup>

در این مرحله، که مرحله نهایی است، پردازش‌های شامل پردازش‌های نحوی<sup>۶۸</sup> برای صیقل‌دادن جملات تولید شده و تصحیح نهایی آن‌ها، صورت می‌گیرد.

اولین روش‌های ارائه شده در این حوزه، عموماً محدود به کاربردهای خاص بودند. در عموم این روش‌ها، مراحل مورد نیاز برای اجرای فرایند تولید جمله به ترتیب زیر، اجرا می‌شدند.

#### ۱. استخراج لغات متناظر با داده از طریق جدول نگاشت ثابت<sup>۶۹</sup>

#### ۲. استخراج عبارات مناسب زبانی برای تولید جمله

#### ۳. اعمال قوانین دستور زبان برای ساخت جمله

<sup>۶۲</sup>Knowledge Bases

<sup>۶۳</sup>Logical Forms

<sup>۶۴</sup>Framework

<sup>۶۵</sup>Text Planning

<sup>۶۶</sup>Sentence Planning

<sup>۶۷</sup>Linguistic Realisation

<sup>۶۸</sup>Syntactic

<sup>۶۹</sup>Hardcoded Mapping Table

در [۲۵]، روشی برای تولید جملات زبان طبیعی، بیان کننده وضعیت حرکت قطارها در یک ترمینال مسافربری، ارائه داده است. در این پژوهش، اطلاعات مختلف موجود در پایگاه داده، به دو دسته پیام‌های «حرکت از مبدأ» و «رسیدن به مقصد» تقسیم می‌شوند و اطلاعات زمانی و شماره هر قطار در هر پیام، استخراج می‌شود. سپس با استفاده از یک جدول از پیش تعیین شده ثابت<sup>۷۰</sup>، هر پیام به یک مجموعه از لغات، نگاشت می‌شود. در ادامه، عبارات زبانی مناسب جهت تولید جملات استخراج شده و با اعمال کلیشه‌های ثابت<sup>۷۱</sup>، جملات نهایی تولید می‌شوند.

#### ۴-۳ روش نزدیک‌ترین همسایه

یکی از پرکاربردترین روش‌ها در این زمینه، استفاده از روش نزدیک‌ترین همسایه است. در این روش، با استفاده از بردار ویژگی‌های به دست آمده از داده‌های مورد تفسیر، و استخراج بردار ویژگی‌های متناظر از جملات موجود در پایگاه داده، نزدیک‌ترین جمله به بردار ویژگی حاصل، به عنوان جمله توصیف‌کننده انتخاب شده و اعلام می‌شود. پژوهش‌های زیادی با استفاده از روش نزدیک‌ترین همسایه، پردازش‌های مختلفی بر روی داده‌های متنی انجام داده‌اند. به عنوان مثال در پژوهش [۲۶] با استفاده از این روش، مدلی جهت تشخیص صحت یا عدم صحت یک جمله به لحاظ دستور زبانی، ارائه شده است. در این پژوهش با استفاده از دو معیار فاصله اقلیدسی و معیار فاصله ویرایش<sup>۷۲</sup>، ارائه شده است. معیار فاصله ویرایش در اصل، میزان هزینه حذف، درج یا تغییر کاراکترها را در یک دنباله کاراکتری برای رسیدن به یک جمله جدید محاسبه می‌کند. مطابق با گزارش این پژوهش، بالاترین دقیقت به دست آمده از این مدل برای تشخیص صحت یک جمله به لحاظ دستور زبانی، برابر با ۵۵٪ بوده است.

فعالیت‌های متعدد دیگری نیز با استفاده از این روش، سعی در شرح تصویر داشته‌اند. برای مثال در پژوهش [۳]، ایده اصلی در تولید شرح بر تصاویر، استفاده از نزدیک‌ترین جمله به تصویر است. روش‌های مختلف و متعدد محاسبه شباهت در این پژوهش مورد بحث و بررسی قرار گرفته‌اند که به اختصار به بیان آن‌ها خواهیم پرداخت. فرض می‌شود یک مجموعه  $S_{cond}$  شامل تمام جملات موجود در مجموعه‌داده که هر کدام شرحی بر یک تصویر هستند و یک مجموعه  $I_{cand}$  شامل تمام جملات موجود در مجموعه‌داده که هر کدام مربوط به یکی از جملات هستند، وجود دارند.

در این پژوهش، دو هدف به طور کلی مورد نظر قرار گرفته است (در هر دو تعریف، تابع  $(i, s)^f$  میزان شباهت جمله  $s$  و تصویر  $i$  را محاسبه می‌کند).

۱. پیدا کردن بهترین جمله توصیف‌کننده یک تصویر

در این مرحله، به دنبال یافتن  $s^*$  به گونه‌ای هستیم که مقدار  $(i, s^*)^f$  به ازای تصویر ورودی  $i$  بیشینه شود.

۲. پیدا کردن بهترین تصویر توصیف‌کننده یک جمله

در این مرحله به دنبال یافتن  $i^*$  به گونه‌ای هستیم که مقدار  $(s, i^*)^f$  به ازای جمله ورودی  $s$  بیشینه شود.

ایده اصلی در این پژوهش، این است که با ورود یک تصویر جدید  $i$  به سیستم، ابتدا نزدیک‌ترین تصویر موجود در مجموعه‌داده را نسبت به این تصویر پیدا کرده  $(i^{NN})$  و سپس نزدیک‌ترین جمله به تصویر بازیابی شده  $(s^{NN})$  را به

<sup>۷۰</sup> Hardcoded Mapping Table

<sup>۷۱</sup> Fixed Templates

<sup>۷۲</sup> Edit Distance

عنوان جمله خروجی انتخاب می‌کنیم. برای پیدا کردن تصویر مرتبط با یک جمله نیز به همین منوال عمل می‌شود.

شکل ۱۸ نتایج اختصاص تصاویر و جملات را توسط این روش نمایش می‌دهد. در این شکل، نتایج بصری به طور کیفی و براساس میزان خطای جمله تولیدی، به چهار دسته تقسیم شده‌اند.

... describes the image without any errors (score = 4)	The selected caption ...	... describes the image with minor errors (score = 3)	... is somewhat related to the image (score = 2)	... is unrelated to the image (score = 1)
				
<i>A girl wearing a yellow shirt and sunglasses smiles.</i>		<i>A Miami basketball player dribbles by an Arizona State player.</i>	<i>A group of people walking a city street in warm weather.</i>	<i>A boy jumps into the blue pool water.</i>
				
<i>A man riding a motor bike kicks up dirt.</i>	<i>Dogs pulling a sled in a sled race.</i>	<i>Two little girls practice martial arts.</i>	<i>A snowboarder in the air over a snowy mountain.</i>	<i>A child jumping on a tennis court.</i>
				
			<i>A boy in a blue life jacket jumps into the water.</i>	<i>A black dog with a purple collar running.</i>

شکل ۱۸: نتایج کیفی اختصاص جملات و تصاویر به یکدیگر با استفاده از روش نزدیکترین همسایه [۳]

به علاوه، جدول ۱ نتایج ارزیابی روش نزدیکترین همسایه را در اختصاص جملات و تصاویر به یکدیگر بر اساس سه معیار نمایش می‌دهد. معیار اول، میانگین امتیازی که افراد خبره به ۱۰۰۰ جملات تولید شده برای هر عکس داده‌اند را نمایش می‌دهد. این امتیازها اعداد صحیح بین ۱ تا ۴ را شامل می‌شوند و بیشینه ممکن برای این معیار، ۴ است. امتیاز بالاتر نشان‌دهنده مناسب‌تر بودن جملات تولید شده هستند. معیارهای ROUGE و BLUE معیارهایی هیتند که در پژوهش‌های ترجمه ماشین به عنوان محکی برای میزان خوب بودن ترجمه تولید شده، استفاده می‌شوند. مقادیر بالاتر در این معیارها، مناسب بودن عملکرد را نمایش می‌دهد.

جدول ۱: نتایج استفاده از روش نزدیکترین همسایه در اختصاص جملات و تصاویر [۳]

ROUGE	BLUE	امتیاز افراد خبره
۰.۱۱	۰.۳۵	۱.۵۷

به عنوان یکی دیگر از روش‌های مورد استفاده در تولید خودکار شرح بر تصاویر، می‌توان به پژوهش ارائه شده در [۴] اشاره کرد. در این پژوهش، تصاویر موجود در مجموعه‌داده، هر کدام با تعدادی عبارت زبانی که توسط کاربران انسانی نوشته شده‌اند، توصیف شده‌اند. در این پژوهش، هدف اصلی این است که با استخراج شبیه‌ترین عبارات موجود در مجموعه‌داده به تصویر ورودی و با کنارهم قرار دادن این عبارات و ساخت جمله با استفاده از چارچوب کاری مورد استفاده در پژوهش‌های تولید زبان طبیعی، جمله متناسب، تولید و نمایش داده شود. علاوه بر استفاده از روش نزدیکترین همسایه برای انتخاب بهترین عبارات زبانی توصیف‌کننده تصویر، می‌توان از

استفاده هوشمندانه از مرحله «برنامه‌ریزی محتوا»<sup>۷۳</sup> در این پژوهش، به عنوان یکی از نقاط قوت آن، یاد کرد. این کار باعث می‌شود علاوه بر تولید جملات سازگار با یکدیگر، از تولید جملات تکراری در یک شرح بر یک تصویر، خودداری شود که از نقاط قوت این روش است. این مرحله با استفاده از روش‌های بهینه‌سازی انجام می‌شود. روش برنامه‌سازی خطی صحیح<sup>۷۴</sup>، به عنوان چارچوب کاری در این مرحله مورد استفاده قرار گرفته است.

در این پژوهش، ۸۹ کلاس از اجسام و ۲۶ کلاس صحنه برای تشخیص محتوا انتخاب شده و تصاویر ورودی، با استفاده از آشکارکننده‌های فلزنسوالب، به این دسته‌ها اختصاص داده می‌شوند. این کار، تخمین مناسبی از محتوای جملات ارائه می‌دهد. همین‌طور با استفاده از روش برچسب‌گذاری نقش کلمات در جمله<sup>۷۵</sup>، محتوای عبارات زبانی متناظر با جملات، به طریق مشابه، دسته‌بندی می‌شود.

چهار دسته از عبارات برای هر تصویر ورودی، به این روش، استخراج می‌شود.

#### ۱. عبارات اسمی

جستجوی عبارات اسمی موجود در مجموعه‌داده با استفاده از ویژگی بافت و رنگ به عنوان معیارهای محاسبه شbahت.

مانند:

”The brown cow”

#### ۲. عبارات فعلی

استفاده از معیارهای مشابه با عبارات اسمی در بین عبارات فعلی موجود در مجموعه‌داده

مانند:

”boy running”

#### ۳. عبارات اضافی نواحی و اجسام

جستجوی عبارات با استفاده از معیارهای شباهت رنگ، بافت، هیستوگرام گرادیان<sup>۷۶</sup> و همین‌طور با درنظر گرفتن ویژگی‌های هندسی اجسام

مانند:

”in the sky” , ”on the road”

#### ۴. عبارات اضافی صحنه<sup>۷۸</sup>

جستجوی عبارات با استفاده از نتیجه دسته‌بندی صحنه با معیار  $L^2$

مانند:

<sup>۷۳</sup>Content Planning

<sup>۷۴</sup>Integer Linear Programming (ILP)

<sup>۷۵</sup>Part of Speech Tagging (POS Tagging)

<sup>۷۶</sup>Region/Stuff Prepositional Phrases

<sup>۷۷</sup>Histogram of Gradient (HOG)

<sup>۷۸</sup>Scene Prepositional Phrases

”at the market” , ”on hot summer day” , ”in Sweden”

هر جمله شامل یک عبارت اسمی، بیان‌کننده مفعول جمله و یک یا چند نمونه از عبارت دیگر بیان‌کننده مفهوم تصویر هستند. چهار نوع عملیات انتزاعی زیر برای ساخت جمله، در نظر گرفته شده است. هر یک از اعمال زیر با استفاده از روش برنامه‌سازی خطی صحیح، انجام می‌شوند.

۱. انتخاب مجموعه مفعولی که نیاز به توصیف دارند (هر مفعول برای یک جمله)

۲. بازچینی و ترتیب‌دهی به مفعول‌ها

۳. انتخاب مجموعه عبارات مورد نیاز برای هر جمله

۴. بازچینی و ترتیب‌دهی عبارات در هر جمله

در شکل ۱۹، نتایج عملکرد الگوریتم را در مقایسه با جملات تولید شده توسط انسان، مشاهده می‌نمایید. مواردی که با عبارت ILP مشخص شده‌اند، خروجی‌های روش برنامه‌سازی خطی صحیح و مواردی که با عبارت HMM نمایش داده شده‌اند، جملات تولید شده توسط انسان را نمایش می‌دهند.



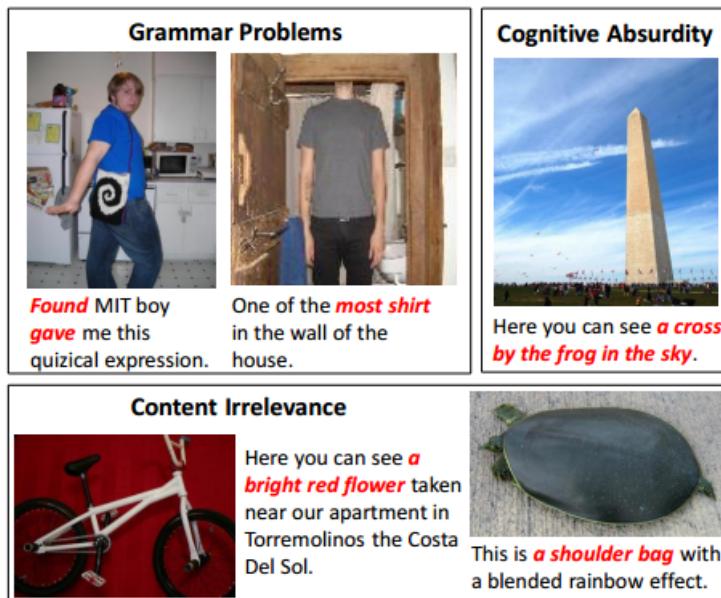
شکل ۱۹: نتایج برنامه‌سازی خطی صحیح در مقایسه با جملات تولید شده انسان [۴]

به علاوه، در شکل ۲۰، مواردی را مشاهده می‌نمایید که در آن‌ها، خروجی الگوریتم نسبت به جملات تولید شده توسط انسان، برتری دارد.



شکل ۲۰: مواردی از خروجی برنامه‌سازی خطی صحیح که نسبت به جملات انسان، برتری دارد [۴].

با وجود این که الگوریتم در برخی موارد کارایی بهتری از خود نشان داده است، جملاتی نیز وجود دارند که به لحاظ دستور زبان، عدم شناخت صحیح یا ناسازگاری محتوایی دچار مشکل شده‌اند. شکل ۲۱، نمونه‌هایی از این دست را نمایش می‌دهد.



شکل ۲۱: نتایج برنامه‌سازی خطی صحیح که به لحاظهای مختلف دچار مشکل شده‌اند [۴].

### ۵-۳ استفاده از قالب‌های آماده زبانی

یکی دیگر از روش‌هایی که برای تولید جمله متناظر با یک تصویر مورد استفاده قرار می‌گیرد، روش استفاده از قالب‌های آماده زبانی است. پژوهش‌هایی که از این روش برای تولید جملات زبانی استفاده کرده‌اند، غالباً پژوهش‌های وظیفه‌محور<sup>۷۹</sup> هستند. حملات تولید شده در این پژوهش‌ها عموماً به طور هدفمند برای پاسخ دادن به موارد معینی تولید می‌شوند و قابلیت تعمیم‌پذیری کمتری نسبت به روش‌های دیگر دارند.

به عنوان مثال در پژوهش [۵]، یک روش مبتنی بر استفاده از قالب‌های آماده زبانی برای تولید خودکار شرح بر تصاویر ارائه شده است. جملاتی که در این پژوهش تولید می‌شوند، باید قادر به مشخص کردن اطلاعات زیر برای هر تصویر باشند:

۱. اجسامی که در تصویر مشاهده می‌شوند

۲. ویژگی‌های اجسام شامل رنگ، اندازه و موارد مشابه

۳. فاعل

۴. فعل

۵. حروف اضافه

<sup>۷۹</sup>Task-Based

در این پژوهش، هدف اصلی این است که پس از بررسی تصویر با استفاده از روش‌های موجود در حاشیه‌نویسی تصویر<sup>۸۰</sup>، که در آن‌ها هر تصویر با یک یا چند برچسب زبانی حاشیه نویسی می‌شود، بتوان جمله‌های بیان‌کننده موارد فوق را در تصویر با استفاده از برچسب‌های تولید شده، تولید نمود. در این پژوهش از مجموعه‌داده PASCAL<sup>۸۱</sup> استفاده شده است که شامل ۱۰۰۰ تصویر و برای هر تصویر، ۵ جمله تولید شده توسط انسان است. ابتدا با استفاده از یک ابزار برچسب‌زنی نقش کلمات در جملات، تمام کلمات موجود در جملات مجموعه‌داده، برچسب زده می‌شوند. سپس برای هر تصویر، ۲ مفعول از بین ۵ جمله متناظر آن و برای هر مفعول، یک ویژگی استخراج می‌شود. به علاوه با استفاده از نرم‌افزار WordNet، مترادف‌های هر یک از کلمات استخراج شده، تا ۳ سطح، یافت می‌شوند.

در مرحله بعدی برای یافتن فاعل جمله، کافیست تعداد دفعاتی را که هر یک از ۲ مفعول استخراج شده، در بین ۵ جمله موجود، در نقش فاعل بوده‌اند شمرده و کلمه‌ای را که بیشترین تعداد تکرار به عنوان فاعل را داشته است، به عنوان فاعل جمله در نظر بگیریم. با مشخص شدن فاعل و مفعول جمله، کافیست فعل مورد نظر را با شمارش تعداد دفعات تکرار افعال مختلف در جملاتی که فاعل و مفعولشان برابر با مورد در حال بررسی است، فعل با بیشترین تکرار را انتخاب کنیم. اگر چنین فعلی یافت نشد، از فعل در قالب آماده استفاده نخواهد شد.

موارد مورد نیاز دیگر نیز به همین ترتیب استخراج می‌شوند. در انتهای، با توجه به پیدا شدن یا نشدن فعل و همین‌طور نوع فعل استخراج شده، از یکی از قالب‌های زیر برای ساخت جمله استفاده می‌شود.

۱. فعل اصلی استخراج شده است:

(معرف ۱ - ویژگی ۱ - فاعل) - فعل - حرف اضافه - (معرف ۲ - ویژگی ۲ - مفعول)

۲. فعل گذرا استخراج شده است:

(معرف ۱ - ویژگی ۱ - فاعل) - فعل - (معرف ۲ - ویژگی ۲ - مفعول)

۳. فعلی یافت نشده است:

(معرف ۱ - ویژگی ۱ - فاعل) - حرف اضافه - (معرف ۲ - ویژگی ۲ - مفعول)

شکل ۲۲ مواردی از جملات تولید شده برای تصاویر موجود در مجموعه‌داده را نمایش می‌دهد. نمونه‌هایی که در این شکل مشاهده می‌شود، نمونه‌هایی هستند که به لحاظ معنایی صحیح بوده و با تصویر مربوطه سازگاری دارند.

<sup>۸۰</sup>Image Annotation

<sup>۸۱</sup><http://vision.cs.uiuc.edu/pascal-sentences/>

green, white, car, saloon	loaded, chair, shelf, plastic	train, green, yellow, building	boat, beach, water, person, sandy	white, airport, engine, building, aeroplane
A green car is parked in front of a white saloon.	A plastic chair is sitting next to a loaded shelf.	A yellow train is passing a green building.	A person is trailing the water. There is a boat in a sandy beach.	A white building is standing on an airport. There is the engine of an aeroplane.

شکل ۲۲: نمونه‌های صحیح از جملات تولید شده توسط قالب‌های آماده زبانی [۵]

به علاوه در شکل ۲۳، نمونه‌هایی از خروجی الگوریتم را در حالتی که جملات تولید شده به لحاظ شناخت صحیح فاعل، ویژگی، فعل، حروف اضافه و همین‌طور تکرار کلمات در جمله دچار مشکل شده‌اند، نمایش می‌دهد.

Incorrect subject	Incorrect attr-obj	Incorrect verb	Incorrect prep	Object Repetition
A white table is sitting on a buddha statue.	A brown person is milking a young cow.	A young person is cutting a cake.	A wooden table is sitting on a wooden chair.	There is a patio set and the wet patio in the picture.

شکل ۲۳: نمونه‌های اشتباه تولید شده توسط قالب‌های آماده زبانی [۵]

علاوه بر این، جدول ۲، نتایج معیار BLUE را در حالات مختلف (استفاده از ترکیبات چندتایی کلمات<sup>۸۲</sup>) نمایش می‌دهد. همان‌طور که مشاهده می‌شود، مقادیر به دست آمده از این معیارها، قابل قبول بودن دقت جملات تولید شده توسط این روش را نمایش می‌دهند. در ستون‌هایی از جدول که از حرف s استفاده شده، تطابق بین کلمات هم‌معنی نیز در نظر گرفته شده است در بقیه ستون‌ها، کلمات دقیق با هم مقایسه شده‌اند. همان‌طور که مشخص است، استفاده از کلمات هم‌معنی نتایج بهتری را ارائه داده است.

جدول ۲: نتایج معیارهای BLUE و ROUGE در حالات مختلف [۵]

ROUGE-1-s	ROUGE-1	BLUE-3-s	BLUE-3	BLUE-2-s	BLUE-2	BLUE-1-s	BLUE-1
۰.۶۰	۰.۵۵	۰.۴۲	۰.۳۵	۰.۶۱	۰.۵۵	۰.۷۹	۰.۷۴

## ۶-۳ روش‌های مبتنی بر شبکه‌های عصبی بازگشتی

اخیراً، استفاده از شبکه‌های عصبی بازگشتی برای تولید جمله، توجه تعداد زیادی از پژوهش‌گران را به خود جلب کرده است. شبکه‌های عصبی بازگشتی، ضمن اثبات قدرت خود در پیش‌بینی سری‌های زمانی، در کاربردهای متعدد و متنوعی مورد استفاده قرار می‌گیرند. از جمله این کاربردها می‌توان به تولید تصاویر، تولید متن، تولید برنامه، تولید موسیقی و مواردی از این دست اشاره نمود. با توجه به ظرفیت بالا و توان یادگیری بالای این مدل‌ها،

<sup>۸۲</sup>n-grams

استفاده از آن‌ها در تولید جملات زبان طبیعی مرتبط با مفهوم یا مفاهیمی خاص، نظر بسیاری را به خود جلب کرده است.

شبکه‌های عصبی بازگشتی در اوخر دهه ۱۹۹۰، ارائه شدند و فعالیت‌های محدودی بر روی مدل‌های کوچکی از آن‌ها مانند مدل Elman انجام شد. به دلیل زمان بالای مورد نیاز برای آموزش این نوع از شبکه‌های عصبی، تا سال ۲۰۱۱، عموم فعالیت‌ها در این زمینه، محدود به استفاده از مدل‌های کوچک از این شبکه‌ها بودند. در سال ۲۰۱۱، آقای هینتون<sup>۸۳</sup> و همکارانش در پژوهش [۶]، با بهره‌گیری از پیشرفت‌های جدید در بهینه‌سازی روش بدون هسین<sup>۸۴</sup> و اعمال این روش‌ها به فرایند تولید جمله در سطح حروف<sup>۸۵</sup> قادر به آموزش یک شبکه عصبی بازگشتی در ۵ روز شدند و نتایج بهتری نسبت به مدل‌های ارائه شده تا آن زمان، حاصل کردند.

در پژوهش [۶]، علاوه بر ارائه یک روش برای آموزش شبکه‌های عصبی بازگشتی عمیق بر مبنای روش بهینه‌سازی بدون هسین، نشان داده شده است که شبکه‌های عصبی بازگشتی استاندارد، عملکرد خوبی در تولید جملات در سطح حروف از خود نشان نمی‌دهند. برای حل این مساله، نوع خاصی از این شبکه‌ها موسوم به شبکه‌های عصبی بازگشتی عمیق ضربی<sup>۸۶</sup> ارائه شده‌اند.

نکات جالبی که در نتایج شبکه عصبی بازگشتی ارائه شده در این پژوهش به چشم می‌خورد، عبارتند از:

۱. تولید ساختارهای زبانی سطح بالا<sup>۸۷</sup>

۲. پشتیبانی از دایره لغات بسیار وسیع

۳. یادگیری تعداد قابل توجهی از دستورات زبانی

۴. تولید تعداد زیادی اسم محتمل که در بین لغات مجموعه‌داده وجود ندارند

۵. توانایی باز و بسته کردن صحیح پرانتزها و نقل قول‌ها در فواصل طولانی بیش از ۳۰ حرف

یک شبکه عصبی بازگشتی را که دنباله زمانی  $(x_1 \dots x_T)$  را به عنوان ورودی گرفته و دنباله  $(h_1 \dots h_T)$  و دنباله  $(o_1 \dots o_T)$  را به ترتیب به عنوان حالات مخفی و خروجی، تولید می‌کند می‌توان مطابق با رابطه ۲۰ بیان کرد. در این رابطه،  $W_{hx}$  ماتریس وزن‌های لایه ورودی به لایه مخفی،  $W_{hh}$  ماتریس وزن‌های لایه مخفی به لایه مخفی (وزن‌های بازگشتی) و  $W_{oh}$  ماتریس وزن‌های لایه مخفی به لایه خروجی است. بردارهای  $b$ ، بردارهای بایاس هستند و مقدار  $W_{hh}h_{t-1}$  در نقطه  $t = 1$  با یک مقدار اولیه جایگزین می‌شود.

$$h_t = \tanh(W_{hx}x_t + W_{hh}h_{t-1} + b_h) \\ o_t = W_{oh}h_t + b_o \quad (۲۰)$$

<sup>۸۳</sup>Hinton

<sup>۸۴</sup>Hessian Free (HF)

<sup>۸۵</sup>Character Level Sentence Generation

<sup>۸۶</sup>Multiplicative Recurrent Neural Networks (MRNN)

<sup>۸۷</sup>High level linguistic structures

از آنجا که مشتقات رابطه ۲۰ قابل محاسبه هستند، به راحتی می‌توان رابطه به روزرسانی وزن‌ها را بر اساس روش نزول در امتداد گرادیان، برای شبکه عصبی بازگشتی، محاسبه کرد. با این وجود، به دلیل رابطه بین پارامترها و دینامیک شبکه، روش نزول در امتداد گرادیان کارایی خوبی در آموزش شبکه برای پیش‌بینی دنباله‌های زمانی در بازه‌های زمانی بزرگ، ارائه نمی‌دهد. به همین دلیل، تا سال ۲۰۱۱ و ارائه روش بهینه‌سازی بدون هسین در [۶] برای آموزش شبکه عصبی بازگشتی، پژوهش‌های زیادی در این زمینه انجام نمی‌شد.

برای رفع مشکل روش پسانشان خطای<sup>۸۸</sup> در آموزش شبکه‌های عصبی بازگشتی، سه روش زیر پیشنهاد شده است:

#### ۱. استفاده از شبکه‌های عصبی حالت پژواک<sup>۸۹</sup>

در این شبکه‌ها، فقط وزن‌های پیش‌رو<sup>۹۰</sup> به روزرسانی می‌شوند. مقادیر اولیه برای وزن‌های بازگشتی باید به طور مناسب انتخاب شوند.

#### ۲. شبکه‌های حافظه کوتاه‌مدت بلند<sup>۹۱</sup>

در این شبکه‌ها، گره‌هایی برای به خاطر سپاری ورودی‌های قدیمی تعبیه می‌شود که باعث افزایش قدرت این شبکه‌ها در پیش‌بینی بلندمدت دنباله‌های زمانی می‌شود.

#### ۳. استفاده از بهینه‌سازی بدون هسین در آموزش شبکه‌های عصبی بازگشتی ضربی<sup>۹۲</sup> این مدل، تاکنون توانسته کارایی بهتری از هر دو مدل قبلی از خود نشان دهد [۶].

در ادامه به بررسی شبکه عصبی ارائه شده در پژوهش [۶] می‌پردازیم.

### ۱-۶-۳ شبکه عصبی بازگشتی ضربی[۹]

ایده اصلی در طراحی این شبکه این است که به جای استفاده و آموزش یک ماتریس یکسان از وزن‌های لایه مخفی به ازای تمام ورودی‌های ممکن، به ازای هر ورودی ممکن، یک ماتریس وزن داشته باشیم. به عبارت دیگر، هر حرکتی که به عنوان ورودی به شبکه وارد شد، مشخص کننده ماتریسی باشد که به عنوان وزن‌های لایه مخفی در آموزش شبکه شرکت می‌کند. با این ایده، رابطه (۲۰) به شکل رابطه (۲۱) تبدیل می‌شود.

$$\begin{aligned} h_t &= \tanh(W_{hx}x_t + W_{hh}^{(x_t)}h_{t-1} + b_h) \\ o_t &= W_{oh}h_t + b_o \end{aligned} \quad (21)$$

همان‌طور که ملاحظه می‌شود، تنها تفاوت رابطه (۲۰) و رابطه (۲۱) در بالا نویس پارامتر  $W_{hh}$  است. پارامتر جدید را می‌توان به شکل رابطه (۲۲) تعریف کرد که در آن،  $x_t^{(m)}$  مولفه  $m$  از ورودی و  $W_{hh}^{(m)}$  است. لازم به

<sup>۸۸</sup>Backpropagation

<sup>۸۹</sup>Echo State Network

<sup>۹۰</sup>Feed forward

<sup>۹۱</sup>Long Short Term Memory (LSTM)

ذکر است در رابطه (۲۲) برای تعریف  $W_{hh}^{(x_t)}$  از تنسور<sup>۹۲</sup> استفاده شده است. در این میان،  $M$  تعداد ابعاد ورودی و ماتریس‌های  $W_{hh}^{(1)}$  تا  $W_{hh}^{(M)}$  ماتریس‌های مربوط به هریک از ابعاد ورودی هستند.

$$W_{hh}^{(x_t)} = \sum_{m=1}^M x_t^{(m)} W_{hh}^{(m)} \quad (22)$$

مدل ارائه شده از آنجا که به طور کامل عمومی است، در حالاتی که ابعاد ورودی و تعداد گره‌های مخفی زیاد باشد، دچار مشکل می‌شود زیرا حجم زیادی حافظه نیاز خواهد داشت. برای حل این مشکل، سعی در فاکتورگیری از پارامتر  $W_{hh}^{(x_t)}$  خواهیم داشت.

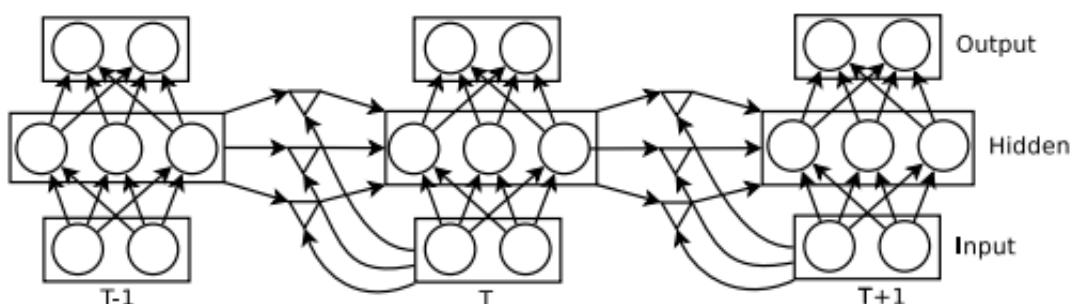
با تعریف سه ماتریس  $W_{hf}$ ،  $W_{fx}$  و  $W_{fh}$  تغییر رابطه (۲۲) به شکل رابطه (۲۳)، می‌توان این فاکتورگیری را تعریف کرد.

$$W_{hh}^{(x_t)} = W_{hf} \cdot \text{diag}(W_{fx} x_t) \cdot W_{fh} \quad (23)$$

در صورتی که ابعاد ماتریس  $W_{fx} x_t$  که آن را  $F$  خواهیم نامید، بزرگ باشد، رابطه (۲۳) به همان اندازه رابطه (۲۲) پیچیده خواهد بود اما در صورتی که این ابعاد کوچک باشد، رابطه (۲۳) نسبت به رابطه (۲۲) به مراتب ساده‌تر خواهد بود. با جای‌گذاری رابطه (۲۳) در رابطه (۲۱)، شبکه عصبی بازگشتی ضربی به دست خواهد آمد. رابطه نهایی را می‌توان به شکل رابطه (۲۴) نمایش داد.

$$\begin{aligned} f_t &= \text{diag}(W_{fx} x_t) \cdot W_{fh} h_{t-1} \\ h_t &= \tanh(W_{hf} f_t + W_{hx} x_t) \\ o_t &= W_{oh} h_t + b_o \end{aligned} \quad (24)$$

شکل ۲۴ طرحواره‌ای از مدل کلی شبکه عصبی بازگشتی ضربی را که توسط رابطه (۲۴) مدل شده است، نمایش می‌دهد. علامت مثلث که در شکل نشان داده شده است، نمایش‌دهنده یک گیت<sup>۹۳</sup> را نمایش می‌دهد که در آن، وزن لایه مخفی با استفاده از ورودی، انتخاب می‌شود.



شکل ۲۴: طرحواره شبکه عصبی بازگشتی ضربی [۶]

<sup>۹۲</sup>Tensor

<sup>۹۳</sup>Gate

نکته‌ای که هم‌چنان مبهم باقی می‌ماند، نحوه محاسبه وزن‌های مخفی است. رابطه (۲۵) نحوه محاسبه  $W_{hh}^{(x_t)ij}$  را نمایش می‌دهد. در این حاصل ضرب، اگر به عنوان مثال پارامتر  $W_{hhif}$  خیلی کوچک باشد و  $W_{hhfj}$  خیلی بزرگ باشد یا برعکس، حاصل مشتق، برای وزن‌های بسیار کوچک، بسیار بزرگ می‌شود و برای وزن‌های بسیار بزرگ، بسیار کوچک. این باعث می‌شود روش نزول در امتداد گرادیان، به حالت پایداری نرسد.

$$W_{hh}^{(x_t)ij} = \sum_f W_{hhif} W_{hhfx^t} W_{hhfj} \quad (25)$$

به دلیل همین عدم پایداری ایجاد شده، حاصل از رابطه (۲۵)، روش نزول در امتداد گرادیان، گزینه مناسبی برای آموزش این شبکه نخواهد بود. لذا لزوم استفاده از روش‌های مرتبه دوم مانند روش بدون هسین، مشهود می‌شود.

### ۲-۶-۳ تولید جمله با مفهوم مشخص

استفاده از شبکه عصبی بازگشتی ضربی منجر به فراهم‌سازی بستری مناسب جهت تولید جمله در سطح حروف می‌شود. با این وجود برای تولید خودکار شرح بر تصاویر، نیازمند آن هستیم که محتوای جملات را به طور مشخص و از پیش تعیین شده داشته باشیم. به همین دلیل نیاز به ارائه روش که طی آن بتوانیم معنا و محتوای جملات تولید شده توسط شبکه عصبی بازگشتی را کنترل کنیم، مشهود می‌شود.

در پژوهش [۷] روش جدیدی برای تولید شرح خودکار بر تصاویر ارائه شده است که در مرحله تولید جمله، از نوع خاصی از شبکه‌های عصبی بازگشتی موسوم به شبکه عصبی بازگشتی دوطرفه<sup>۹۴</sup>، استفاده شده است. در این روش، ابتدا از یک شبکه عصبی کانولوشنی عمیق بر روی نواحی استخراج شده از تصاویر برای استخراج ویژگی و درک صحنه استفاده شده است. از سوی دیگر، با اعمال یک شبکه عصبی بازگشتی دوطرفه بر جملات و ارائه یک تابع هدف ساختارمند، روشی برای همترازسازی جمله و اطلاعات بصری نهفته در تصویر ارائه شده است.

شکل ۲۵، نمونه‌ای از همترازسازی ارائه شده در این پژوهش برای یک تصویر را نشان می‌دهد.

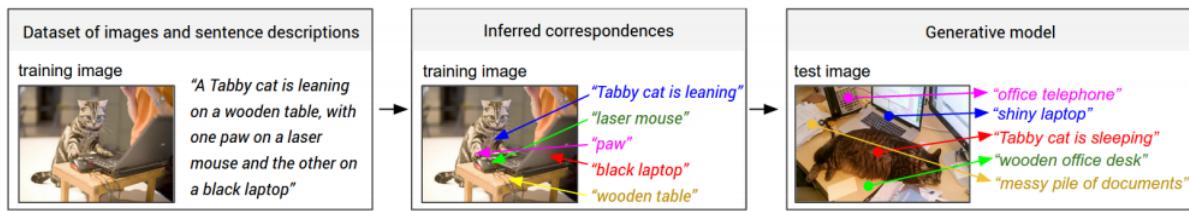


شکل ۲۵: همترازسازی تصویر و جمله [۷]

<sup>۹۴</sup>Bidirectional RNN

در مدل همترازسازی ارائه شده در این پژوهش، فرض بر این است که یک مجموعه‌داده شامل تعداد زیادی تصویر و جملات متناظر با هر تصویر وجود دارد. همین‌طور فرض دیگری وجود دارد مبنی بر این که بخش‌های مختلف هر جمله، به نواحی خاصی از تصویر اشاره می‌کنند که موقعیت این نواحی مجھول است. از طرف دیگر اشاره این بخش‌ها به نواحی مرتبط خود در تصاویر، در بین تمام مجموعه‌داده، تکرار می‌شود. به عنوان مثال، عبارات زبانی شامل کلمه «توب» در تمام تصاویر موجود در مجموعه‌داده، به نواحی از تصویر اشاره می‌کنند که دارای ویژگی‌های «توب» هستند.

در شکل ۲۶ ارتباط بین بخش‌های مختلف یک جمله و نواحی متفاوت از تصویر را مشاهده می‌نمایید. همان‌طور که در شکل مشاهده می‌شود، ابتدا برای تصاویر موجود در مجموعه‌داده و شرح متناظر با هر یک از این تصاویر، ارتباطات بین عبارات مختلف از جملات و نواحی تصاویر استخراج و یادگرفته می‌شود. در ادامه، با ورود یک تصویر جدید و براساس ارتباطات یادگرفته شده، شرح جدید برای تصویر تولید می‌شود.



شکل ۲۶: ارتباط بین نواحی مختلف یک تصویر و عبارات جمله<sup>[۷]</sup>

برای تبدیل تصویر به فضای ویژگی، مطابق با آن‌چه در فصل درک صحنه ذکر شد، ابتدا با استفاده از روش شبکه‌های عصبی کانولوشنی ناحیه‌ای، ۱۹ ناحیه از تصویر استخراج شده و برای ۲۰ تصویر موجود، با اعمال یک بهینه‌سازی و تخمین پارامتر و اعمال یک شبکه عصبی، بردار ویژگی استخراج می‌شود. پس از استخراج بردار ویژگی از نواحی تصویر، نیازمند آن هستیم که بتوانیم از عبارات مختلف جمله، بردار ویژگی هماندازه با بردار ویژگی حاصل از تصویر، استخراج کنیم. برای این کار، در این پژوهش از شبکه عصبی بازگشتی دوطرفه استفاده شده است.

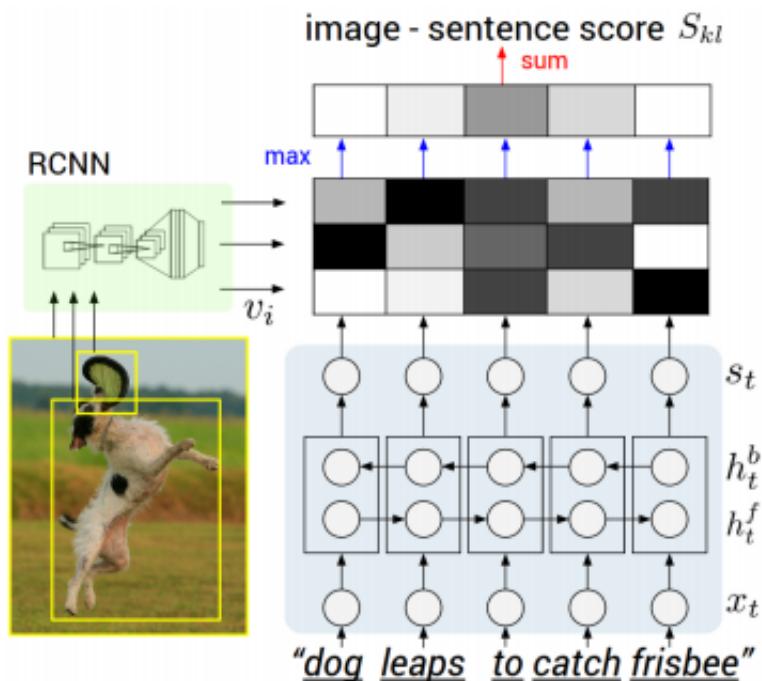
این شبکه عصبی، یک دنباله از  $N$  کلمه را به عنوان ورودی دریافت کرده و هر یک را به یک بردار در فضای  $h$  بعدی، که  $h$  اندازه بردار ویژگی حاصل از نواحی تصویر است، نگاشت می‌کند. رابطه (۲۶)، رابطه مربوط به پارامترهای این شبکه عصبی را نمایش می‌دهد. در این رابطه،  $I$  یک بردار ستونی اندیکاتور<sup>۹۵</sup> است که در اندیس کلمه  $t$ ام خود یک و در بقیه اندیس‌ها صفر دارد.  $W_w$  یک ماتریس وزن ثابت برای هر کلمه  $w$  است که برای جلوگیری از بیش‌برازش بر داده‌ها، مورد استفاده قرار می‌گیرد.

<sup>۹۵</sup>Indicator

$$\begin{aligned}
x_t &= W_w I_t \\
e_t &= f(W_e x_t + b_e) \\
h_t^f &= f(e_t + W_f h_{t-1}^f + b_f) \\
h_t^b &= f(e_t + W_b h_{t-1}^b + b_b) \\
s_t &= f(W_d(h_t^f + h_t^b) + b_d)
\end{aligned} \tag{26}$$

در این شبکه، دو جریان داده وجود دارد. جریان اول، جریان داده بین گره‌های مخفی شبکه از چپ به راست و دیگری جریان داده بین نودهای مخفی شبکه از راست به چپ است که به ترتیب با  $h_t^f$  و  $h_t^b$  نمایش داده می‌شوند. بردار نهایی  $s_t$ ، بردار حاصل از نگاشت کلمات به فضای ویژگی‌ها است که با استفاده از خود کلمه و محتوای مورد استفاده در اطراف کلمه در جمله، تولید می‌شود.

شکل ۲۷ طرح‌واره‌ای از معماری این شبکه را نمایش می‌دهد. همان‌طور که در این شکل مشخص است، در لایه مخفی این شبکه، دو جریان داده، یکی از راست به چپ و دیگری از چپ به راست برای محاسبه تاثیر کلمات اطراف کلمه جاری بر نگاشت کلمه به فضای ویژگی‌ها، وجود دارد.



شکل ۲۷: طرح‌واره شبکه عصبی بازگشتی دوطرفه [۷]

در ادامه با بهره‌گیری از روش نگاشت دوطرفه تصاویر و جملات که در بخش درک صحنه ارائه شد، توابع هم‌ترازسازی و تابع هدف ارائه شده را مورد استفاده قرار داده و با استفاده از روش یادگیری چند نمونه‌ای، اقدام به یادگیری انتساب‌های بین نواحی مختلف تصاویر و عبارات مختلف زبانی می‌شود. با استفاده از این روش، می‌توان برای هر ناحیه از تصویر، کلمات مناسب را تعیین کرد. اما برای تولید خودکار شرح

بر تصاویر، نیاز به تولید عبارات زبانی وجود دارد. برای حل این مشکل، با در نظر گرفتن رابطه ضرب داخلی بین بردارهای ویژگی حاصل از نواحی تصویر و عبارات زبانی یک جمله به عنوان معیار شباht، روشی ارائه شده است که بتوان برای هر ناحیه از تصویر، عبارت زبانی مناسبی تولید کرد.

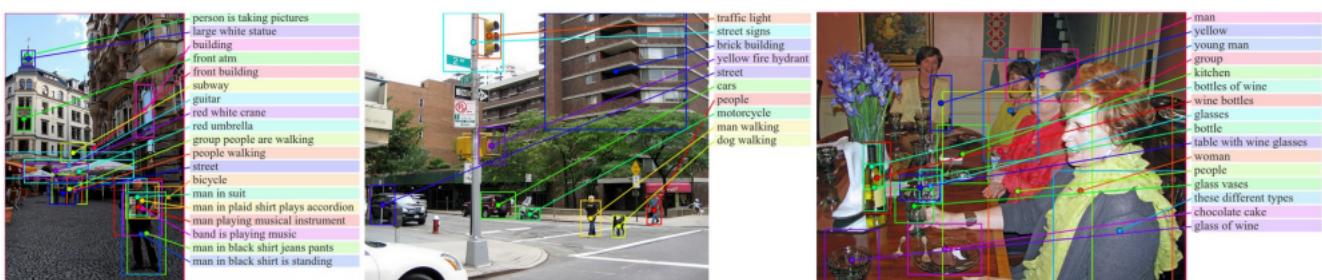
در این روش، با تعریف یک تابع انرژی و استفاده از مدل میدان تصادفی مارکف، با بهینه‌سازی تابع انرژی ارائه شده، بهترین هم‌ترازسازی برای هر یک از عبارات موجود محاسبه شده و عبارت با بهترین مقدار، انتخاب می‌شود. رابطه (۲۷)، این تابع انرژی را محاسبه می‌کند.

$$E(a) = \sum_{j=1 \dots N} \Psi_j^U(a_j) + \sum_{j=1 \dots M} \Psi_j^B(a_j, a_{j+1})$$

$$\Psi_j^U(a_j) = \nu_i^T \cdot s^t$$

$$\Psi_j^B(a_j, a_{j+1}) = \beta I(a_j = a_{j+1}) \quad (27)$$

می‌توان در شکل ۲۸ نتایج استفاده از این شبکه و محاسبه میزان شباهت نواحی مختلف تصویر و عبارات مختلف از جملات را مشاهده نمود. این شکل، نتیجه آموزش اختصاص نواحی مختلف تصویر به عبارات زبانی را نمایش می‌دهد.



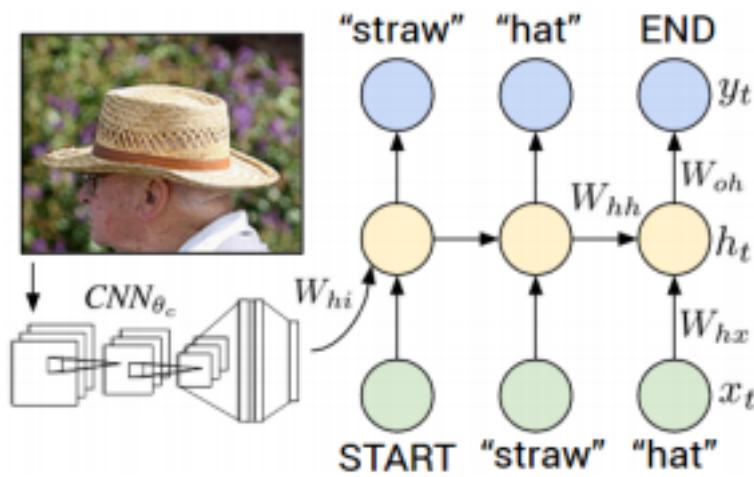
شکل ۲۸: انتساب نواحی مختلف تصویر به عبارات زبانی [۷]

تا این مرحله، با ورود یک تصویر، یا یک ناحیه از تصویر، عبارات زبانی متناظر، استخراج شده‌اند. هدف اصلی در این پژوهش تولید جمله برای هر تصویر است. بنابراین نیاز داریم تا با استفاده از مدل‌های ارائه شده برای تولید جمله، این کار را انجام دهیم. در پژوهش‌های زیادی، استفاده از شبکه‌های عصبی بازگشتی برای پیش‌بینی و محاسبه توزیع احتمال کلمه بعدی در یک جمله با در نظر داشتن کلمات قبلی و محتوای جمله، ارائه شده است. در این پژوهش با اعمال تغییرات کوچکی، از همین روش‌ها استفاده می‌شود.

رابطه ارائه شده برای شبکه عصبی بازگشتی که این کار را انجام می‌دهد، مطابق با رابطه (۲۸) است. در این رابطه،  $CNN_{\theta_c}(IMAGE)$  بردار حاصل از اعمال آخرین لایه یک شبکه عصبی کانولوشنی بر تصویر را نشان می‌دهد و بقیه پارامترها، همگی قابل آموزش هستند. بردار  $y_t$  بردار نماینده توزیع احتمالاتی تمام کلمات با در نظر گرفتن کلمات قبلی و محتوای هر ناحیه است که اندازه آن برابر است با تعداد تمام کلمات موجود در لغتنامه به علاوه یک نشانه خاص به عنوان «اتمام جمله».

$$\begin{aligned}
 b_\nu &= W_{hi}[CNN_{\theta_c}(IMAGE)] \\
 h_t &= f(W_{hx}x_t + W_{hh}h_{t-1} + b_h + I(t=1) \cdot b_\nu) \\
 y_t &= softmax(W_{oh}h_t + b_o)
 \end{aligned} \tag{28}$$

شکل ۲۹، طرحواره‌ای از شبکه عصبی بازگشتی مالتی‌مدال<sup>۹۶</sup> ارائه شده در این پژوهش برای تولید جمله را نمایش می‌دهد.



شکل ۲۹: طرحواره شبکه عصبی بازگشتی ارائه شده برای تولید جمله [۷]

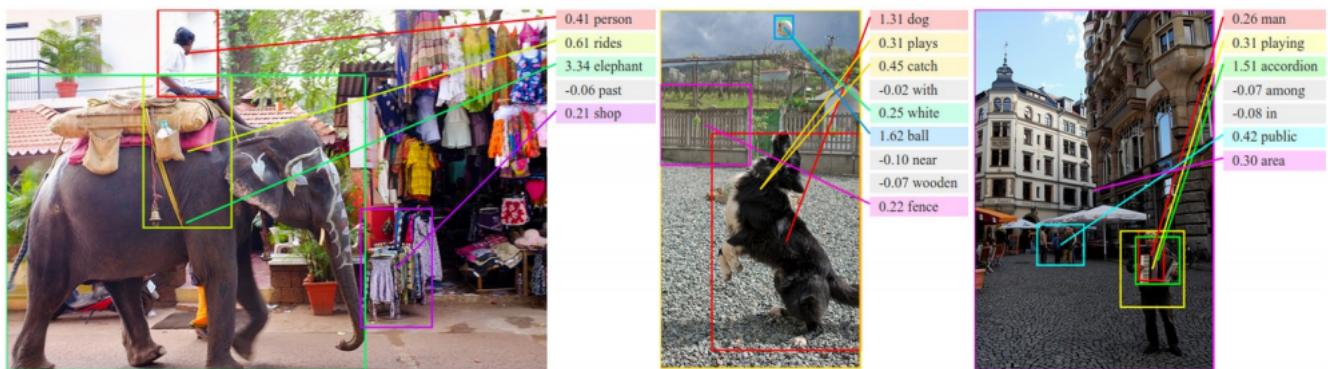
جدول ۳ نتایج معیار BLUE را برای روش ارائه شده بر روی سه مجموعه‌داده مختلف، نمایش می‌دهد. همان‌طور که در این جدول مشخص است، نتایج به دست آمده از این روش در مقایسه با دو روش دیگر، نتایج به نسبت بهتری بوده است.

جدول ۳: نتایج معیار BLUE برای روش ارائه شده در [۷] در مقایسه با دو روش دیگر

B-4**	B-3**	B-2**	B-1**	B-3*	B-2*	B-1*	B-3	B-2	B-1	نام روش
۱۰۰	۱۶.۶	۲۸.۱	۴۵.۰	—	—	—	—	—	—	نرديک‌ترین‌همسايه
—	—	—	—	۲۰	۲۴	۵۵	۲۳	۲۸	۵۸	روش [۲۷]
۲۳.۰	۳۲.۱	۴۵.۰	۶۲.۵	۲۴.۰	۳۶.۹	۵۷.۳	۲۴.۵	۳۸.۳	۵۷.۹	روش [۷]

شکل ۳۰ نتایج رتبه‌بندی جملات و عبارات برای هر ناحیه از تصویر را نمایش می‌دهد.

<sup>۹۶</sup>Multimodal



شکل ۳۰. نتایج رتبه‌بندی عبارات زبانی برای نواحی تصویر [۷]

به علاوه، در شکل ۳۱ نتایج تولید شرح برای تصاویر، توسط شبکه عصبی بازگشته ارائه شده در این پژوهش، به تصویر کشیده شده است.



شکل ۳۱. نتایج تولید جمله برای تصاویر در [۷]

### ۷-۳ جمع‌بندی

چالش تولید جمله یکی از قدیمی‌ترین و پویاترین حوزه‌های فعالیتی و پژوهشی در هوش مصنوعی است که از اواسط قدن بیستم، توجه پژوهش‌گران بسیاری را به خود جلب کرده است. روش‌های مختلفی برای حل این مساله ارائه شده‌اند. از جمله این روش‌ها می‌توان به موارد زیر اشاره کرد:

#### ۱. تولید زبان طبیعی

در این دسته از روش‌ها که از اواخر دهه بیستم تا کنون مورد استفاده قرار می‌گیرند، با طی فرایند در یک چارچوب کلی، سعی در تولید جملات مناسب دارند. این دسته از روش‌ها عموماً برای تفسیر خودکار داده‌هایی که برای کاربران انسانی غیر قابل تفسیر هستند یا تفسیر دشواری دارند، به کار می‌روند. در این روش‌ها ابتدا با استفاده از ویژگی‌های مختلفی که در داده‌های ماشینی (داده‌های قابل تفسیر برای ماشین) کلمات مناسب انتخاب شده و سپس با استفاده از کلمات منتخب، عبارات زبانی (با جایگشت دادن کلمات و حذف عبارات غیر محتمل) تولید می‌شوند. سپس با اعمال قواعد دستور زبان و چینش عبارات زبانی در کاره‌هم، جملات نهایی تولید می‌شوند.

## ۲. نزدیکترین همسایه

در این دسته از روش‌ها سعی می‌شود با ورود یک تصویر و نگاشت آن به فضای ویژگی‌ها، جمله‌ای از میان تمام جملات موجود در مجموعه‌داده انتخاب شود که بیشترین مشابهت با بردار ویژگی تصویر را دارد. بزرگترین مشکل در این روش‌ها انتخاب معیار مناسب برای محاسبه فاصله بین یک جمله و بردار ویژگی حاصل از تصویر است. در این روش، علاوه بر این که نیاز به وجود مجموعه‌داده وسیع و پوشای وجود دارد، ممکن است جمله نهایی، در انتهای گویا و بیان‌کننده تمام جوانب تصویر نباشد و یا حتی با تصویر ورودی سازگاری نداشته باشد.

برای حل این مشکل، سعی شد به جای استخراج نزدیکترین جمله به تصویر موجود، مشابه‌ترین عبارات زبانی را با شکستن جملات موجود به عبارات سازنده، انتخاب کرده و با بهره‌گیری از روش تولید زبان طبیعی و یا روش‌های دیگر، چینش مناسبی از این عبارات را که در قالب یک یا چند جمله بیان شوند، تولید و به عنوان شرح بر تصویر، نمایش داد.

## ۳. استفاده از قالب‌های زبانی آماده

با وجود فعالیت‌های گوناگون در این زمینه و استفاده از روش‌های مختلف، هم‌چنان تضمین صحت جمله خروجی، کار دشواری است. به همین دلیل، سعی شد با ارائه یک یا چند قالب زبانی آماده و از پیش تعیین شده برای جملات، مانند قالب‌های جملات خبری، صحت جملات نهایی را تضمین کرد. در این دسته از روش‌ها، ویژگی‌های مختلفی از تصویر استخراج می‌شود که هریک از این ویژگی‌ها یا همه آن‌ها در کنار هم قادر هستند نقش‌هایی مانند « فعل »، « فاعل »، « مفعول » و موارد مشابه را در جمله متناظر با تصویر مشخص کنند. با استخراج کلمات مناسب و شناخت نقش آن‌ها در جمله و جای‌گذاری هر یک از این کلمات در مکان مناسب نقشی خود در قالب از پیش تعیین شده، جمله متناظر با هر تصویر استخراج می‌شود.

## ۴. استفاده از شبکه‌های عصبی بازگشتی

اگر چه استفاده از قالب‌های آماده و از پیش تعیین شده، تا حدی مشکلات موجود را حل می‌کند اما هم‌چنان چالش بزرگ‌تری حل نشده باقی مانده است. تولید جملات جدید، استفاده از کلمات و عبارات جدید و ابتکاری به طوری که علاوه بر تضمین رعایت دستور زبان، بتوان معنای جمله را نیز متضمن شد، چالش بزرگی است که در این مسیر کماکان وجود دارد.

استفاده از شبکه‌های عصبی بازگشتی یکی از بهترین راه‌کارهای موجود برای حل این مشکل و رویارویی با این چالش هستند. استفاده از این شبکه‌ها در اواخر قرن بیستم در بین پژوهش‌گران رواج پیدا کرد تا جایی که ناپایداری الگوریتم پساننتشار خطا در آموزش این شبکه، راه را برای پژوهش‌های بعدی بست. پس از ارائه یک روش مناسب برای بهینه‌سازی بدون هسین در سال ۲۰۱۰، روشی برای آموزش یک شبکه عصبی بازگشتی موسوم به شبکه عصبی بازگشتی ضربی بر مبنای بهینه‌سازهای بدون هسین ارائه شد و نتایج آن به طور چشم‌گیری از روش‌های موجود بیشتر بود.

ارائه شبکه عصبی بازگشتی ضربی، نقطه عطفی در مسیر علم در راستای حل چالش تولید جمله به حساب می‌آید. از حدود سال ۲۰۱۱ به بعد، استفاده از شبکه‌های عصبی بازگشتی برای تولید جمله به پویاترین و پرفعالیت‌ترین حوزه در مسائل مربوط به تولید جمله، به حساب می‌آید.

خانم لی و همکارانش در سال ۲۰۱۵، در پژوهش [۷]، با استفاده از شبکه‌های عصبی کانولوشنی عمیق و دو نوع از شبکه‌های عصبی بازگشتی موسوم به شبکه‌های عصبی بازگشتی مالتی‌مدال و شبکه‌های عصبی بازگشتی دوطرفه، روش مناسبی برای تولید خودکار شرح بر تصاویر ارائه داده است.

در این پژوهش، ابتدا با بهره‌گیری از روش شبکه عصبی کانولوشنی ناحیه‌ای، نواحی از تصویر که شامل تصویر اجسام است، استخراج شده و با استفاده از یک شبکه عصبی کریشفسکی، بردار ویژگی برای هر ناحیه محاسبه می‌شود. سپس با بهره‌گیری از یک شبکه عصبی بازگشتی دوطرفه، عبارات مختلف از جمله استخراج و بردارهای ویژگی برای هر عبارت محاسبه می‌شود. سپس با استفاده از یک تابع هدف و مدل میدان تصادفی مارکف، همترازسازی بین نواحی و عبارات زبانی صورت گرفته و مدل آموزش داده می‌شود.

در ادامه با تخمین بهینه پارامترهای موجود و با استفاده از شبکه عصبی بازگشتی مالتی‌مدال، توزیع احتمال بهترین کلمه بعدی در یک جمله با داشتن کلمات قبلی و محتوای حاصل از بردار ویژگی محاسبه شده روی نواحی تصویر، محاسبه شده و بهترین کلمه بعدی تولید می‌شود. این کار تا جایی ادامه می‌یابد که شبکه، نشانه مخصوص پایان جمله را تولید کند.

## ۴ فصل چهارم

تولید شرح متناظر صحنه با استفاده از  
یادگیری عمیق

## ۱-۴ تولید شرح بر تصاویر با استفاده از یادگیری عمیق

## ۵ فصل پنجم

تولید شرح متناظر صحنه با استفاده از  
روش‌های مبتنی بر نقطه توجه

## ۱-۵ تولید شرح بر تصاویر با استفاده از روش‌های مبتنی بر نقطه توجه

ایده اصلی روش‌های مبتنی بر نقطه توجه از پژوهش‌های موجود در زمینه ترجمه ماشینی گرفته شده است. این دسته از پژوهش‌ها مدلی ارائه می‌دهند که با استفاده از آن بتوان هر کلمه از جملات تولیدی را با تمرکز بر یک یا بخشی از کلمات موجود در جمله مبدأ، تولید کرد. به طور مشابه، در حوزه تولید خودکار شرح بر تصاویر، از این دسته از پژوهش‌ها به منظور حصول مدلی استفاده می‌شود که قادر باشد هر یک از کلمات موجود در جمله را با استفاده از بخشی از تصویر ورودی، تولید نماید.

در این فصل، ابتدا ایده اصلی ترجمه مبتنی بر نقطه توجه را در حوزه ترجمه ماشینی ارائه خواهیم کرد و سپس کاربردهای این ایده را در حوزه تولید خودکار شرح بر تصاویر مورد بررسی قرار می‌دهیم.

## ۲-۵ روش‌های مبتنی بر نقطه توجه در حوزه ترجمه ماشینی

همان‌طور که گفته شد، تمام روش‌های قبلی را می‌توان به دو مرحله زیر تقسیم کرد.

۱. نگاشت نمونه‌ها از فضای تصاویر به فضای ویژگی‌ها

۲. نگاشت نمونه‌ها از فضای ویژگی‌ها به فضای جملات

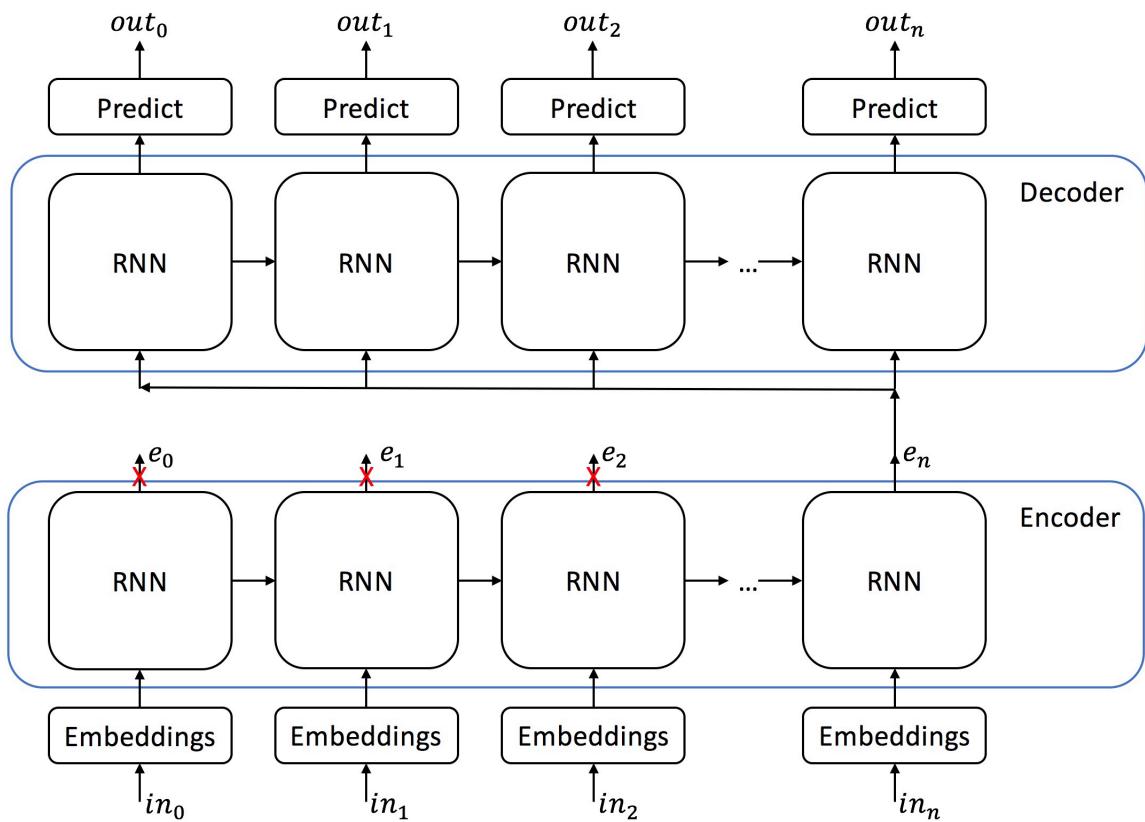
در حوزه ترجمه ماشینی، به تابع نگاشت مرحله اول، انکودر<sup>۹۷</sup> و به تابع نگاشت مرحله دوم، دیکودر<sup>۹۸</sup> گفته می‌شود. در این بخش، ما از این عبارات برای ارجاع به مراحل اول و دوم الگوریتم استفاده می‌نماییم. چارچوب کاری انکودر-دیکودر، در تعداد زیادی از پژوهش‌های حوزه ترجمه ماشینی به عنوان چارچوب کاری اصلی مورد استفاده قرار گرفته است. تمام روش‌های قبلی که در فصول قبل ذکر شد نیز از همین چارچوب کاری به عنوان چارچوب اصلی بهره برده‌اند. به عنوان مثال در روش‌های مبتنی بر یادگیری عمیق برای تولید خودکار شرح بر تصاویر از شبکه‌های عصبی کانولوشنی به طور معمول به عنوان انکودر و از شبکه‌های عصبی بازگشتی به عنوان دیکودر استفاده می‌شود.

شكل ؟؟ نشان‌دهنده ساختار کلی چارچوب کاری انکودر-دیکودر است.

در ادامه به بررسی بخش‌های مختلف این چارچوب کاری می‌پردازیم و سپس ایده اصلی روش‌های مبتنی بر نقطه توجه را که توسط آقای بنجیو در پژوهش [۲۸] در سال ۲۰۱۴ ارائه شده است، مورد بررسی قرار خواهیم داد.

<sup>۹۷</sup>Encoder

<sup>۹۸</sup>Decoder



شکل ۳۲: ساختار کلی چارچوب کاری انکودر-دیکودر

## ۱-۲-۵ انکودر

انکودر در این چارچوب کاری، با گرفتن یک جمله به عنوان ورودی، بردار ویژگی متناظر جمله مبدا را تولید می‌کند. جمله ورودی با دنباله‌ای از کلمات مدل می‌شود. همین‌طور هر کلمه را با یک بردار  $n$  بعدی، که  $n$  تعداد کلمات موجود در دیکشنری است، مدل می‌شود. به این ترتیب، هر جمله ورودی، یک بردار با طول متغیر است که هر مولفه آن خودش برداری به ابعاد  $n$  است. از طرفی بردار خروجی، که همان بردار ویژگی‌ها است، یک بردار با طول ثابت و قراردادی خواهد بود.

عموماً در کاربردهای ترجمه ماشینی در هر دو بخش انکودر و دیکودر از شبکه‌های عصبی بازگشتی استفاده می‌شود. در شبکه‌های عصبی بازگشتی، خروجی هر مرحله تابعی از ورودی آن مرحله و حالت شبکه در مرحله جاری است. با فرض این‌که  $h_t$  حالت شبکه در زمان  $t$  را نمایش دهد می‌توان رابطه تولید خروجی توسط شبکه عصبی بازگشتی را مطابق با (۲۹) تعریف نمود.

$$h_t = f(X_t, h_{t-1})$$

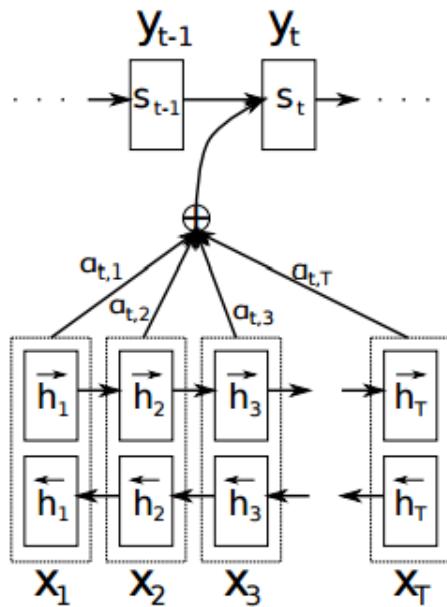
$$C = q(h_1, \dots, h_L) \quad (29)$$

به طور معمول از شبکه LSTM به عنوان تابع  $f$  استفاده می‌شود و همین‌طور به جای استفاده از تابع  $q$  حالتنهایی شبکه به عنوان بردار ویژگی مورد استفاده قرار می‌گیرد [۲۸].

دیکودر به منظور نگاشت فضای ویژگی‌ها به فضای جملات مورد استفاده قرار می‌گیرد. خروجی انکودر، ورودی دیکودر است. با این فرض، ورودی دیکودر یک بردار ویژگی با طول ثابت است و خروجی آن که یک جمله به زبان مقصد است، همانند جمله مبدا، یک بردار با طول متفاوت شامل بردارهای بازنمایی کلمات است. دیکودر در اصل در هر مرحله، به دنبال یافتن کلمه‌ای است که با داشتن کلمات تولید شده قبلی و بردار ویژگی موجود، محتمل‌ترین کلمه نسبت به بقیه کلمات موجود در دیکشنری باشد.تابع احتمال مربوطه را می‌توان به فرم (۳۰) تعریف نمود.

$$p(y_t|C, y_1, y_2, \dots, y_{t-1}) = g(y_t, s_t, C) \quad (30)$$

در رابطه (۳۰)،  $C$  نشان‌دهنده بردار ویژگی،  $y_i$  نشان‌دهنده لغت  $i$ ام تولید شده از زبان مقصد و بردار  $s_t$  نشان‌دهنده حالت شبکه بازگشتی مورد استفاده به عنوان دیکودر است.  
شکل ۳۳ ساختار کلی دیکودر را نمایش می‌دهد.



شکل ۳۳: ساختار دیکودر مورد استفاده در چارچوب کاری [۲۸]

### ۳-۲-۵ ایده اصلی استفاده از نقطه توجه

همان‌طور که بیان شد، در چارچوب کاری انکودر-دیکودر، ابتدا جمله ورودی که شامل تعداد نامعلوم کلمه است به یک بردار با طول متفاوت مدل می‌شود. بردار تولید شده توسط یک انکودر به یک بردار با طول ثابت، که همان بردار ویژگی‌ها است، نگاشت شده و در نهایت بردار ویژگی تولید شده توسط یک انکودر به یک بردار با طول متفاوت که نماینده جمله زبان مقصد است، نگاشت می‌شود.  
فرآیند مذکور یک محدودیت جدی دارد و آن این است که انکودر باید بتواند تمام اطلاعات مورد نیاز برای تولید

جمله را در یک بردار با طول ثابت بگنجاند و دیکودر باید بتواند تمام اطلاعات مورد نیاز خود را از همین بردار با موجود با طول ثابت، استخراج کند. این محدودیت باعث می‌شود قدرت کد کردن اطلاعات در بردار ویژگی کاهش یابد. برای حل این مشکل از ایده نقاط توجه استفاده می‌نماییم.

در این دسته از روش‌ها به جای این‌که انکودر فقط یک بردار ویژگی تولید کند، بردارهای ویژگی مختلفی ایجاد می‌کند که هر بردار با مرکز بر روی یک یا بخشی از جمله مبدا تولید شده است. به این ترتیب، هر بردار تولید شده شامل اطلاعات معنایی یک یا بخشی از جمله مبدا می‌باشد. به این طریق، دیکودر می‌تواند با انتخاب بین بردارهای معنایی تولید شده در هر مرحله، کلمه تولیدی را با مرکز بر روی معنای یک کلمه و کلمات مجاور آن در جمله مبدا، تولید کند.

در ادامه به بررسی تغییراتی که باید در انکودر و دیکودر اتفاق بیفتد تا بتوان به جای یک بردار ویژگی مجموعه‌ای از بردارهای ویژگی با مرکز محلی ایجاد نمود و از آن‌ها برای تولید جمله استفاده نمود را مورد بررسی قرار می‌دهیم. برای سهولت فهم تغییرات، ابتدا تغییرات دیکودر را مطرح نموده و سپس به بررسی تغییرات انکودر خواهیم پرداخت.

#### ۴-۲-۵ دیکودر در روش مبتنی بر نقطه توجه

فرض می‌کنیم به جای تنها یک بردار ویژگی،  $L$  بردار ویژگی از ورودی استخراج شده باشد. آن‌ها را در یک ماتریس به شکل  $C = [c_1, \dots, c_L]^T$  بازنمایی می‌نماییم. فرض می‌کنیم بردار ویژگی  $c_i$  به دنباله حاشیه‌نویسی‌های<sup>۹۹</sup>  $h = [h_1, \dots, h_L]^T$  وابسته است. حاشیه‌نویسی  $h_i$  خود یک متغیر تصادفی به شکل برداری است که دارای دو ویژگی بسیار مهم می‌باشد.

۱. حاوی اطلاعات استخراج شده از تمام جمله است

۲. مرکز استخراج اطلاعات بر روی کلمه نام و کلمات اطراف آن بوده است.

با تعریف این دو ویژگی، حاشیه‌نویسی‌ها را می‌توان همان بردار ویژگی جمله تصور کرد با این شرط که علاوه بر این که معنای کل جمله را کد کرده‌اند، مرکز بیشتری بر معنای کلمه نام و کلمات مجاور آن دارند. به عبارت بهتر هر حاشیه‌نویسی علاوه بر این‌که معنای کلی جمله را کد می‌کند، حاوی معنای محلی مربوط به کلمات هم هست. با تعریف حاشیه‌نویسی به شکل فوق و با تکیه بر فرض‌های انجام شده، می‌توانیم مدل احتمالاتی ارائه شده را به شکل (۳۱) تغییر دهیم.

$$p(y_i|y_1, \dots, y_{i-1}, X) = g(y_{i-1}, S_i, c_i) \quad (31)$$

که در آن:

$$c_i = \sum_{j=1}^L \alpha_{ij} h_j \quad (32)$$

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^L \exp(e_{ik})} \quad (33)$$

$$e_{ij} = f(s_{i-1}, h_j) \quad (34)$$

---

<sup>۹۹</sup>Annotation

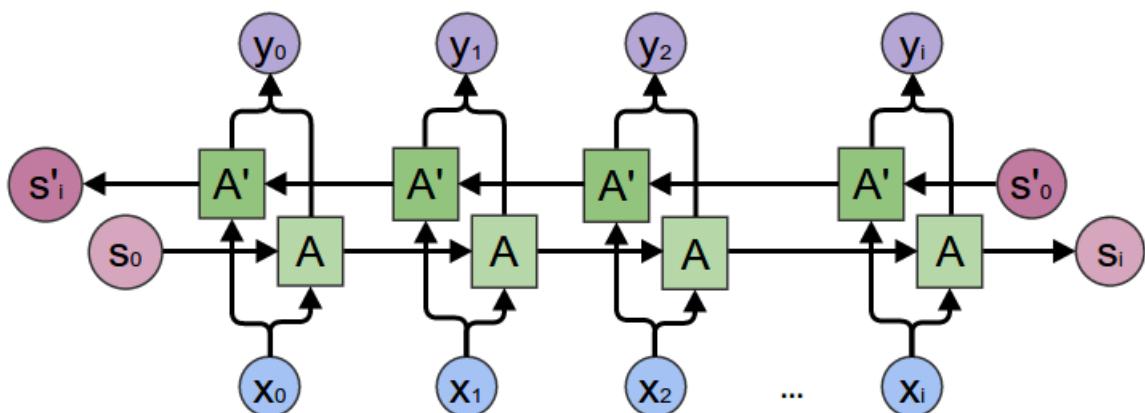
متغیر تصادفی  $e_{ij}$  که در رابطه (۳۴) تعریف شده است نمایان‌گر میزان شباهت کلمه  $z$ ام در جمله خروجی به کلمه  $z$ ام در جمله ورودی است. وظیفه این متغیر، هم‌ترازسازی<sup>۱۰۰</sup> ورودی و خروجی است.  $\alpha_{ij}$  یک نرمال‌سازی روی امتیازهای محاسبه شده انجام می‌دهد. از این متغیر نرمال شده به عنوان وزن حاشیه‌نویسی‌ها استفاده می‌شود. مطابق با رابطه (۳۲) بردار ویژگی مورد استفاده برای تولید کلمه در جمله مقصود، از طریق یک میانگین‌گیری بر اساس وزن معنایی کلمات تولید می‌شود.

در رابطه (۳۴)  $s_{i-1}$  بردار حالت شبکه دیکودر در زمان  $i-1$  و  $f$  یکتابع امتیاز است. تابع امتیاز مورد استفاده در این رابطه را می‌توان با یک شبکه عصبی پیش‌رو<sup>۱۰۱</sup> مدل‌سازی کرد. در صورت استفاده از شبکه عصبی پیش‌رو برای مدل‌سازی تابع شباهت، در صورتی که از هم‌ترازسازی نرم<sup>۱۰۲</sup> استفاده شود، تابع هدف مشتق‌پذیر شده و می‌توانیم از الگوریتم پساننتشار خطاب برای آموزش استفاده نماییم.

## ۵-۲-۵ انکودر در روش مبتنی بر نقطه توجه

برای طراحی انکودر در این بخش، باید مکانیزمی ارائه شود که قادر باشد حاشیه‌نویسی‌های  $h_1$  تا  $h_L$  را طوری تولید کند که دو شرط مطرح شده در بخش قبلی را ارضاء نمایند. به عبارت دیگر باید بردارهای ویژگی‌ای استخراج نماییم که علاوه بر این‌که حاوی معنای کل جمله باشند، هر یک از آن‌ها بر روی معنای یک کلمه و کلمات اطراف آن تمرکز بیشتری نسبت به سایر بردارها داشته باشند تا بتوانیم علاوه بر مدل‌سازی معنای کلی جمله، از معنای محلی کلمات هم استفاده نماییم.

به این منظور از یک شبکه عصبی بازگشتی دوطرفه در مدل‌سازی انکودر استفاده می‌نماییم. شکل ۳۴ ساختار کلی یک شبکه عصبی بازگشتی دوطرفه را نمایش می‌دهد.



شکل ۳۴: ساختار کلی یک شبکه عصبی بازگشتی دوطرفه

همان‌طور که در شکل ۳۴ مشخص است، یک شبکه عصبی بازگشتی دوطرفه شامل دو شبکه پیش‌رو در خلاف جهت یک‌دیگر است. حالت‌های مخفی شبکه پیش‌رو را به راست را با  $\rightarrow h_i$  و حالت‌های مخفی شبکه پیش‌رو را به چپ را با  $\leftarrow h_i$  نمایش می‌دهیم. همان‌طور که در شکل پیداست، خروجی‌های شبکه در این ساختار هم به حالت‌های

<sup>۱۰۰</sup> Alignment

<sup>۱۰۱</sup> Feed Forward Neural Network

<sup>۱۰۲</sup> Soft Alignment

سمت راست و کلمات سمت راست در جمله و هم به حالات و کلمات سمت چپ وابسته هستند. پس همین خروجی‌ها را می‌توان به عنوان حاشیه‌نویسی‌هایی که هر دو ویژگی را دارند مورد استفاده قرار داد.

## ۶ فصل ششم

تولید شرح متناظر صحنه با استفاده از  
روش‌های مبتنی بر حافظه فعال

## ١-٦ تولید جمله

## ۷ فصل هفتم

### جمع‌بندی و نتیجه‌گیری

## ۱-۷ جمع‌بندی و نتیجه‌گیری

در این مستند، گزارش مختصری درباره مساله تولید خودکار شرح بر تصاویر و روش‌های پیشنهادی برای حل چالش‌های موجود در این مسیر را ارائه دادیم. مساله تولید خودکار شرح بر تصاویر، به معنای تولید جملات زبان طبیعی برای هر تصویر است به طوری که این جملات شامل سه شرط زیر باشند:

۱. صحنه، اجسام موجود در تصویر، رابطه مکانی اجسام و اطلاعاتی از این دست که به درک تصویر کمک می‌کنند باید در جملات تولید شده وجود داشته باشند و دقیق و کامل باشند.
۲. جملات تولید شده، خود، باید به لحاظ معنایی، دستور زبان و املایی صحیح بوده و نقصی نداشته باشند.
۳. جملات تولید شده باید با تصاویر مرتبط با خود، سازگاری داشته باشند.

ایده‌های اولیه در این مسیر از پژوهش‌های موجود در زمینه ترجمه ماشین ایجاد شده است که در آن‌ها، ابتدا یک جمله ورودی از یک زبان مبدا، با استفاده از روش‌های مختلفی به یک بردار ویژگی تبدیل می‌شود و سپس در مرحله دوم، بردار ویژگی حاصل، با استفاده از روش‌های خاص دیگری به جملات زبان طبیعی به زبان مقصد، تبدیل می‌شوند. حال اگر به جای جمله از زبان مبدا، یک تصویر به این سامانه وارد شود و با روشی بتوان این تصویر را به همان بردار ویژگی نگاشت کرد، جمله نهایی معادل معنایی تصویر ورودی خواهد شد. با استفاده از این فرایند می‌توان به طور خودکار برای تصاویر، شرح مناسبی به زبان طبیعی تولید کرد.

در این مسیر، دو چالش عمده در پیش رو وجود دارد که باید مرتفع شوند:

۱. چالش درک صحنه  
فرایند استخراج اطلاعات بصری نهفته در تصویر و بازنمایی مناسب این اطلاعات را به گونه‌ای که بتواند برای پردازش‌های بعدی مناسب باشد، فرایند درک صحنه می‌نامند. روش‌های مختلف و متعددی برای حل این چالش، تاکنون مطرح شده‌اند. در این مساله، باید بتوان تصاویر ورودی را به نحوی موثر و مفید به فضای ویژگی‌ها نگاشت کرد به طوری که بازنمایی حاصل، بتواند در مرحله تولید جمله، منجر به تولید جملات معنادار و مناسب شود.

۲. چالش تولید جمله  
تولید جملات به زبان طبیعی که علاوه بر صحت معنایی، دستور زبانی و املایی، قادر به توصیف و تفسیر اطلاعات غیر قابل تفسیر برای کاربران انسانی هستند، از جمله مهم‌ترین و پویاترین حوزه‌های پژوهشی در زمینه هوش مصنوعی است و توجه پژوهش‌گران بسیاری را به خود جلب کرده است. در این مساله، باید

بتوان بردار ویژگی حاصل از تصویر را که در مرحله درک صحنه تولید شده است، به نحوی کارا و موثر به یک جمله در زبان طبیعی نگاشت کرد.

اولین مرحله از فرایند تولید خودکار شرح برای تصاویر، مرحله درک صحنه است. در این مرحله، تصاویر ورودی تحت عملیات مختلفی به فضای معنایی نگاشت می‌شوند. فضای معنایی در اینجا، می‌تواند فضای شامل میدان‌های اطلاعاتی از پیش تعیین شده (مانند فضای سه‌تایی‌های «جسم، رخداد، صحنه») یا فضای بردار ویژگی‌ها باشد. روش‌های مختلفی برای نگاشت تصویر ورودی به فضای معنایی ارائه شده است که به طور کلی می‌توان عموم آن‌ها را به دو بخش تقسیم کرد:

#### ۱. روش‌های مبتنی بر مدل‌های گرافی احتمالاتی

در این روش‌ها با استفاده از مدل‌های استاندارد گرافی احتمالاتی موجود یا با ارائه یک مدل گرافی احتمالاتی، تصویر ورودی به فضای معنایی نگاشت می‌شود. در روش‌های مبتنی بر این مدل‌ها، با ارائه یک توزیع احتمال برای نقاط مختلف در فضای معنایی، محتمل‌ترین نقطه برای تصویر به عنوان نقطه نظر تصویر، انتخاب می‌شود.

##### (آ) مدل میدان تصادفی مارکف

یک نمونه از روش‌های مبتنی بر مدل میدان تصادفی مارکف که برای درک صحنه از آن استفاده شده است، در پژوهش [۱۳] ارائه شده است. درک صحنه در این پژوهش با ارائه یک سه‌تایی «جسم، فعالیت، صحنه» به‌ازای هر تصویر، تعریف شده است. مبتنی بر همین تعریف، یک مدل میدان تصادفی مارکف شامل سه گره که دوبهدو به هم متصل هستند، تعریف شده است. هر یک از گره‌های موجود در این مدل، نماینده یکی از میدان‌های سه‌گانه تعریف شده در فضای معنایی هستند. با تعریف توابع پتانسیل مختلف روی هر گره و توابع پتانسیل مختلف روی هر یال، یک تابع توزیع توام برای تمام متغیرهای تصادفی موجود در مدل ارائه شده است.

با محاسبه مقادیر پتانسیل برای تصاویر مختلف موجود در مجموعه آموزشی و با استفاده از یک ماشین بردار پشتیبان، بردارهای ویژگی شاخص برای هر گره محاسبه می‌شوند. از این بردارهای ویژگی بعداً برای انطباق تصاویر با مقادیر مختلف در هر گره استفاده می‌شود.

در این پژوهش، با یافتن نزدیک‌ترین همسایه‌های یک تصویر بر حسب معیار شباهت با بردارهای ویژگی شاخص و میانگین‌گیری روی مقادیر هر گره، بهترین انطباق تصویر و نقاط فضای معنایی به دست می‌آید. به این ترتیب، برای هر تصویر ورودی، می‌توان نقطه نظر در فضای معنایی را مشخص کرد.

##### (ب) مدل میدان تصادفی شرطی

در پژوهش [۱۲] یک مدل میدان تصادفی شرطی سلسله‌مراتبی برای درک صحنه ارائه شده است که شامل دو سطح انتزاع است. برای گره‌های موجود در هریک از سطوح انتزاع مدل، یک دسته متغیر تصادفی تعریف شده و برای کل مدل سه نوع تابع پتانسیل مختلف معرفی شده است.

اولین دسته از توابع پتانسیل معرفی شده در این بخش، توابع پتانسیل قطعه‌بندی یگانی هستند که به منظور یکپارچه‌سازی نقاط داخل یک قطعه تعریف شده‌اند. توابع پتانسیل دیگری برای انطباق بین متغیرهای تصادفی موجود در بین دو سطح انتزاع تعریف شده‌اند که در صورت مغایرت مقادیر اختصاص داده شده به متغیرهای موجود بین دو سطح، مقدار  $\lambda$  – و در غیر این صورت مقدار صفر دارند. این توابع در شرایطی که مقادیر متغیرهای موجود در دو سطح با هم یکسان نباشد، یک مقدار جریمه به تابع هدف اضافه می‌کنند. آخرین دسته از توابع پتانسیل مورد استفاده، برای انطباق تصویر با دسته تشخیص داده شده اجسام تعریف شده است که توسط فلزنسوالب ارائه شده و به روش دی پیام مشهور است.

(ج) سایر مدل‌های گرافی احتمالی در پژوهش [۲]، یک مدل گرافی احتمالی مولد برای نگاشت تصویر به فضای معنایی ارائه شده است. در این مدل، از دو سطح تصویر استفاده شده است؛ تصویر سطح جسم و تصویر سطح صحنه. برای تصویر سطح صحنه، یک متغیر تصادفی، بیان‌کننده دسته صحنه و برای تصویر سطح جسم دو متغیر تصادفی، بیان‌کننده دسته و شکل جسم، ارائه شده است. روابط بین متغیرهای تصادفی در این پژوهش، براساس نحوه تولید متغیرهای تصادفی و روابط منطقی موجود بین آن‌ها طراحی شده‌اند.

تصویر ورودی در این پژوهش، ابتدا به نواحی کوچک  $10 \times 10$  تقسیم می‌شود و مطابق با روش توضیح داده شده، مقدار توابع پتانسیل مختلف برای هر کدام از متغیرهای تصادفی، در هر ناحیه، محاسبه می‌شود. در این پژوهش، یک تابع احتمال شرطی برای متغیرهای تصادفی ارائه شده است که در مرحله استنتاج، با استفاده از روش تخمین بیشترین احتمال، برچسب‌های هر تصویر مشخص می‌شوند.

۲. روش‌های مبتنی بر استفاده از شبکه‌های عصبی کانولوشنی عمیق  
در این روش‌ها، با ارائه یک شبکه عصبی کانولوشنی عمیق و تعریف کردن تابع هدف برای شبکه، تابع نگاشت تصویر و فضای معنا تشکیل می‌شود. پس از ارائه تابع هدف برای هر شبکه، با بهینه‌سازی آن تابع، پارامترهای موجود در شبکه آموزش داده می‌شوند.

در پژوهش [۲۰]، روشی ارائه شده است که طی آن یک تصویر، به نواحی کوچک‌تر تقسیم می‌شود به طوری که هر ناحیه به وجود آمده، به طور یکپارچه، حاوی یک جسم باشد و هر جسم تنها در یک ناحیه قرار بگیرد. این روش موسوم به روش RCNN است. در این روش، دو ویژگی برای یک ناحیه‌بندی خوب در تصاویر ارائه شده است و پیرو این ویژگی‌ها، روشی برای طرح نواحی پیشنهادی در یک تصویر که دارای این دو ویژگی باشد، ارائه شده است.

ویژگی مطرح شده اول برای ناحیه‌بندی تصاویر این است که، ناحیه‌های ایجاد شده در هر تصویر، می‌توانند در ابعاد مختلف وجود داشته باشند زیرا اجسام موجود در تصاویر، ممکن است اندازه و تعداد متفاوتی داشته باشند. دومین ویژگی برای یک ناحیه‌بندی خوب، این است که معیار انتخاب نواحی نباید برای تمام تصاویر، یکسان در نظر گرفته شود؛ زیرا معیارهای مختلف برای ناحیه‌بندی تصاویر در شرایط مختلف، رفتارهای متفاوتی از خود نشان می‌دهند. بنابراین باید از معیارهای مختلف برای تعیین نواحی استفاده نمود. در این پژوهش، ابتدا تصاویر مطابق با یک معیار اولیه، به مجموعه‌ای از نواحی اولیه تقسیم می‌شوند. سپس

با استفاده از معیارهای مختلف مانند فضاهای رنگی مختلف، معیارهای شباهت مختلف و نقاط اولیه متفاوت، با پیروی از یک روش حریصانه، نواحی کوچکتر که به یکدیگر شبیه‌تر هستند با هم ترکیب شده و نواحی بزرگتر را می‌سازند. نواحی ایجاد شده در این روش، سپس به یک شبکه عصبی کانولوشنی عمیق داده می‌شوند و برای هر ناحیه، یک بردار ویژگی ۴۰۹۶ بعدی ایجاد می‌شود که هر ناحیه با آن بازنمایی شود. در پژوهش [۲۳] با استفاده از روش RCNN و تعریف دوتابع هدف دیگر برای شبکه، روشی ارائه شده است که طی آن بتوان تصاویر و جملات را به طور دوطرفه به یکدیگر نگاشت کرد. توابع هدف تعریف شده در این پژوهش، دوتابع مختلف هستند. اولین تابع هدف، یک تابع هدف سراسری است. این تابع به این منظور تعریف شده است که تصاویر و جملاتی که مطابق با محاسبات شبکه عصبی ارائه شده، بیشترین شباهت را با یکدیگر دارند، در واقعیت هم شبیه‌ترین تصاویر و جملات به یکدیگر باشند. تابع هدف دوم برای این شبکه به این شکل تعریف شده است که نواحی استخراج شده از تصویر و عبارات استخراج شده از جملات که در روش ارائه شده، بیشترین شباهت را به یکدیگر دارند، در واقعیت هم بیشترین شباهت و ارتباط را با یکدیگر داشته باشند.

در این پژوهش، تصاویر ورودی با استفاده از روش RCNN به نواحی مختلف تقسیم شده و ۱۹ ناحیه با بیشترین اطمینان از بین این نواحی انتخاب می‌شود. این ۱۹ ناحیه به همراه خود تصویر به عنوان ۲۰ تصویر مختلف مورد استفاده قرار می‌گیرند. جملات ورودی با استفاده از روشی که در فصل تولید جملات زبان طبیعی توضیح داده خواهد شد، به عبارات مختلف تقسیم می‌شوند و بین هر عبارت استخراج شده و هر یک از ۲۰ تصویر موجود، یک معیار شباهت محاسبه شده و بیشترین شباهت‌ها با هم درنظر گرفته می‌شوند. معیار شباهت مورد استفاده در این روش، ضرب داخلی بین بردارهای ویژگی عبارات و نواحی است. عبارات و نواحی که بیشترین شباهت را با یکدیگر دارند برای تولید جمله به مرحله بعد، ارسال می‌شوند.

چالش تولید جمله یکی از قدیمی‌ترین و پویاترین حوزه‌های فعالیتی و پژوهشی در هوش مصنوعی است که از اواسط قدن بیستم، توجه پژوهش‌گران بسیاری را به خود جلب کرده است. روش‌های مختلفی برای حل این مساله ارائه شده‌اند. از جمله این روش‌ها می‌توان به موارد زیر اشاره کرد:

## ۱. تولید زبان طبیعی

در این دسته از روش‌ها که از اواخر دهه بیستم تا کنون مورد استفاده قرار می‌گیرند، با طی فرایند در یک چارچوب کلی، سعی در تولید جملات مناسب دارند. این دسته از روش‌ها عموماً برای تفسیر خودکار داده‌هایی که برای کاربران انسانی غیر قابل تفسیر هستند یا تفسیر دشواری دارند، به کار می‌روند. در این روش‌ها ابتدا با استفاده از ویژگی‌های مختلفی که در داده‌های ماشینی (داده‌های قابل تفسیر برای ماشین) کلمات مناسب انتخاب شده و سپس با استفاده از کلمات منتخب، عبارات زبانی (با جایگشت دادن کلمات و حذف عبارات غیر متحمل) تولید می‌شوند. سپس با اعمال قواعد دستور زبان و چینش عبارات زبانی در کنارهم، جملات نهایی تولید می‌شوند.

## ۲. نزدیک‌ترین همسایه

در این دسته از روش‌ها سعی می‌شود با ورود یک تصویر و نگاشت آن به فضای ویژگی‌ها، جمله‌ای از میان تمام جملات موجود در مجموعه داده انتخاب شود که بیشترین مشابهت با بردار ویژگی تصویر را دارد. بزرگ‌ترین

مشکل در این روش‌ها انتخاب معیار مناسب برای محاسبه فاصله بین یک جمله و بردار ویژگی حاصل از تصویر است. در این روش، علاوه بر این که نیاز به وجود مجموعه‌داده وسیع و پوشش وجود دارد، ممکن است جمله نهایی، در انتهای گویا و بیان‌کننده تمام جوانب تصویر نباشد و یا حتی با تصویر ورودی سازگاری نداشته باشد.

برای حل این مشکل، سعی شد به جای استخراج نزدیک‌ترین جمله به تصویر موجود، مشابه‌ترین عبارات زبانی را با شکستن جملات موجود به عبارات سازنده، انتخاب کرده و با بهره‌گیری از روش تولید زبان طبیعی و یا روش‌های دیگر، چینش مناسبی از این عبارات را که در قالب یک یا چند جمله بیان شوند، تولید و به عنوان شرح بر تصویر، نمایش داد.

### ۳. استفاده از قالب‌های زبانی آماده

با وجود فعالیت‌های گوناگون در این زمینه و استفاده از روش‌های مختلف، همچنان تضمین صحت جمله خروجی، کار دشواری است. به همین دلیل، سعی شد با ارائه یک یا چند قالب زبانی آماده و از پیش تعیین شده برای جملات، مانند قالب‌های جملات خبری، صحت جملات نهایی را تضمین کرد. در این دسته از روش‌ها، ویژگی‌های مختلفی از تصویر استخراج می‌شود که هریک از این ویژگی‌ها یا همه آن‌ها در کنار هم قادر هستند نقش‌هایی مانند « فعل »، « فاعل »، « مفعول » و موارد مشابه را در جمله متناظر با تصویر مشخص کنند. با استخراج کلمات مناسب و شناخت نقش آن‌ها در جمله و جای‌گذاری هر یک از این کلمات در مکان مناسب نقشی خود در قالب از پیش تعیین شده، جمله متناظر با هر تصویر استخراج می‌شود.

### ۴. استفاده از شبکه‌های عصبی بازگشتی

اگر چه استفاده از قالب‌های آماده و از پیش تعیین شده، تا حدی مشکلات موجود را حل می‌کند اما همچنان چالش بزرگ‌تری حل نشده باقی مانده است. تولید جملات جدید، استفاده از کلمات و عبارات جدید و ابتکاری به طوری که علاوه بر تضمین رعایت دستور زبان، بتوان معنای جمله را نیز متنضم شد، چالش بزرگی است که در این مسیر کماکان وجود دارد.

استفاده از شبکه‌های عصبی بازگشتی یکی از بهترین راه‌کارهای موجود برای حل این مشکل و رویارویی با این چالش هستند. استفاده از این شبکه‌ها در اواخر قرن بیستم در بین پژوهش‌گران رواج پیدا کرد تا جایی که ناپایداری الگوریتم پسانشمار خطا در آموزش این شبکه، راه را برای پژوهش‌های بعدی بست. پس از ارائه یک روش مناسب برای بهینه‌سازی بدون هسین در سال ۲۰۱۰، روشی برای آموزش یک شبکه عصبی بازگشتی موسوم به شبکه عصبی بازگشتی ضربی بر مبنای بهینه‌سازهای بدون هسین ارائه شد و نتایج آن به طور چشم‌گیری از روش‌های موجود بیشتر بود.

ارائه شبکه عصبی بازگشتی ضربی، نقطه عطفی در مسیر علم در راستای حل چالش تولید جمله به حساب می‌آید. از حدود سال ۲۰۱۱ به بعد، استفاده از شبکه‌های عصبی بازگشتی برای تولید جمله به پویاترین و پرفعالیت‌ترین حوزه در مسائل مربوط به تولید جمله، به حساب می‌آید.

خانم لی و همکارانش در سال ۲۰۱۵، در پژوهش [۷]، با استفاده از شبکه‌های عصبی کانولوشنی عمیق و دو نوع از شبکه‌های عصبی بازگشتی موسوم به شبکه‌های عصبی بازگشتی مالتی‌مدال و شبکه‌های عصبی بازگشتی

دوطرفه، روش مناسبی برای تولید خودکار شرح بر تصاویر ارائه داده است.

در این پژوهش، ابتدا با بهره‌گیری از روش شبکه عصبی کانولوشنی ناحیه‌ای، نواحی از تصویر که شامل تصویر اجسام است، استخراج شده و با استفاده از یک شبکه عصبی کریشفسکی، بردار ویژگی برای هر ناحیه محاسبه می‌شود. سپس با بهره‌گیری از یک شبکه عصبی بازگشتی دوطرفه، عبارات مختلف از جمله استخراج و بردارهای ویژگی برای هر عبارت محاسبه می‌شود. سپس با استفاده از یک تابع هدف و مدل میدان تصادفی مارکف، همترازسازی بین نواحی و عبارات زبانی صورت گرفته و مدل آموزش داده می‌شود.

در ادامه با تخمین بهینه پارامترهای موجود و با استفاده از شبکه عصبی بازگشتی مالتی‌مدال، توزیع احتمال بهترین کلمه بعدی در یک جمله با داشتن کلمات قبلی و محتوای حاصل از بردار ویژگی محاسبه شده روی نواحی تصویر، محاسبه شده و بهترین کلمه بعدی تولید می‌شود. این کار تا جایی ادامه می‌یابد که شبکه، نشانه مخصوص پایان جمله را تولید کند.

## مراجع

- [1] Fei-Fei, Li, Iyer, Asha, Koch, Christof, and Perona, Pietro. What do we perceive in a glance of a real-world scene? *Journal of vision*, 7(1):10–10, 2007.
- [2] Li, Li-Jia and Fei-Fei, Li. What, where and who? classifying events by scene and object recognition. in *Computer Vision, 2007. ICCV 2007. IEEE 11th International Conference on*, pp. 1–8. IEEE, 2007.
- [3] Hodosh, Micah, Young, Peter, and Hockenmaier, Julia. Framing image description as a ranking task: Data, models and evaluation metrics. *Journal of Artificial Intelligence Research*, 47:853–899, 2013.
- [4] Kuznetsova, Polina, Ordóñez, Vicente, Berg, Alexander C, Berg, Tamara L, and Choi, Yejin. Collective generation of natural image descriptions. in *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers- Volume 1*, pp. 359–368. Association for Computational Linguistics, 2012.
- [5] Gupta, Ankush and Mannem, Prashanth. From image annotation to image description. in *International Conference on Neural Information Processing*, pp. 196–204. Springer, 2012.
- [6] Sutskever, Ilya, Martens, James, and Hinton, Geoffrey E. Generating text with recurrent neural networks. in *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pp. 1017–1024, 2011.
- [7] Karpathy, Andrej and Fei-Fei, Li. Deep visual-semantic alignments for generating image descriptions. in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3128–3137, 2015.
- [8] Xu, Kelvin, Ba, Jimmy, Kiros, Ryan, Cho, Kyunghyun, Courville, Aaron, Salakhutdinov, Ruslan, Zemel, Richard S, and Bengio, Yoshua. Show, attend and tell: Neural image caption generation with visual attention. *arXiv preprint arXiv:1502.03044*, 2(3):5, 2015.
- [9] Hoiem, Derek, Hays, James, Xiao, Jianxiong, and Khosla, Aditya. Guest editorial: Scene understanding. *International Journal of Computer Vision*, 112(2):131–132, 2015.

- [10] Potter, Mary C. Short-term conceptual memory for pictures. *Journal of experimental psychology: human learning and memory*, 2(5):509, 1976.
- [11] Potter, Mary C, Staub, Adrian, Rado, Janina, and O'Connor, Daniel H. Recognition memory for briefly presented pictures: the time course of rapid forgetting. *Journal of Experimental Psychology: Human Perception and Performance*, 28(5):1163, 2002.
- [12] Fidler, Sanja, Sharma, Abhishek, and Urtasun, Raquel. A sentence is worth a thousand pixels. in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1995–2002, 2013.
- [13] Farhadi, Ali, Hejrati, Mohsen, Sadeghi, Mohammad Amin, Young, Peter, Rashtchian, Cyrus, Hockenmaier, Julia, and Forsyth, David. Every picture tells a story: Generating sentences from images. in *Computer Vision–ECCV 2010*, pp. 15–29. Springer, 2010.
- [14] Felzenszwalb, Pedro, McAllester, David, and Ramanan, Deva. A discriminatively trained, multiscale, deformable part model. in *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pp. 1–8. IEEE, 2008.
- [15] Divvala, Santosh K, Hoiem, Derek, Hays, James H, Efros, Alexei A, and Hebert, Martial. An empirical study of context in object detection. in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*, pp. 1271–1278. IEEE, 2009.
- [16] Lin, Dahua, Fidler, Sanja, and Urtasun, Raquel. Holistic scene understanding for 3d object detection with rgbd cameras. in *The IEEE International Conference on Computer Vision (ICCV)*, December 2013.
- [17] Ladický, Lubor, Sturgess, Paul, Alahari, Karteek, Russell, Chris, and Torr, Philip HS. What, where and how many? combining object detectors and crfs. in *Computer Vision–ECCV 2010*, pp. 424–437. Springer, 2010.
- [18] Ladicky, Lubor, Russell, Chris, Kohli, Pushmeet, and Torr, Philip HS. Graph cut based inference with co-occurrence statistics. in *Computer Vision–ECCV 2010*, pp. 239–253. Springer, 2010.
- [19] Felzenszwalb, Pedro F, Girshick, Ross B, McAllester, David, and Ramanan, Deva. Object detection with discriminatively trained part-based models. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 32(9):1627–1645, 2010.

- [20] Girshick, Ross, Donahue, Jeff, Darrell, Trevor, and Malik, Jitendra. Rich feature hierarchies for accurate object detection and semantic segmentation. in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2014.
- [21] Uijlings, Jasper RR, van de Sande, Koen EA, Gevers, Theo, and Smeulders, Arnold WM. Selective search for object recognition. *International journal of computer vision*, 104(2):154–171, 2013.
- [22] Krizhevsky, Alex, Sutskever, Ilya, and Hinton, Geoffrey E. Imagenet classification with deep convolutional neural networks. in *Advances in neural information processing systems*, pp. 1097–1105, 2012.
- [23] Karpathy, Andrej, Joulin, Armand, and Li, Fei Fei F. Deep fragment embeddings for bidirectional image sentence mapping. in *Advances in neural information processing systems*, pp. 1889–1897, 2014.
- [24] Karpathy, Andrej and Fei-Fei, Li. Deep visual-semantic alignments for generating image descriptions. in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- [25] Reiter, Ehud and Dale, Robert. Building applied natural language generation systems. *Natural Language Engineering*, 3(01):57–87, 1997.
- [26] Lawrence, Steve, Fong, Sandiway, and Giles, C Lee. Natural language grammatical inference: a comparison of recurrent neural networks and machine learning methods. in *International Joint Conference on Artificial Intelligence*, pp. 33–47. Springer, 1995.
- [27] Mao, Junhua, Xu, Wei, Yang, Yi, Wang, Jiang, and Yuille, Alan L. Explain images with multimodal recurrent neural networks. *arXiv preprint arXiv:1410.1090*, 2014.
- [28] Bahdanau, Dzmitry, Cho, Kyunghyun, and Bengio, Yoshua. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.