

## توصیف تصاویر مبتنی بر شبکه عمیق رمزگذار-رمزگشا و سازوکار توجه بر توجه

زهرا فامیل ستاری، حسن ختن لو و الهام علیقارداش

### چکیده

توصیف تصویر یک زمینه تحقیقاتی بین رشته ای در بینایی ماشین و پردازش زبان طبیعی است. بسیاری از روش های پیشنهاد شده برای تولید توصیف تصویر از چارچوب رمزگذار - رمزگشا پیروی کرده اند. به این ترتیب هر کلمه بر اساس ویژگی های تصویر و کلمات تولید شده قبلی تولید می شود. اخیراً سازوکار توجه، که میتواند با ایجاد نقشه فضایی، مناطق مرتبط تصویر با هر کلمه را برجسته کند، به طور گسترده در تحقیقات استفاده شده است. در این مقاله، ما یک روش جدید را پیشنهاد کرده ایم که چارچوب رمزگذار-رمزگشا را با سازوکار توجه و سازوکار توجه بر توجه ادغام کرده است. بخش رمزگذار مدل شامل چند بخش ResNet، Attention-LSTM، Multi Head Attention و Attention on Attention است. از ResNet برای استخراج ویژگی های کلی تصویر استفاده شده است. ایهی Language-LSTM مسئولیت رمزگشایی را بر عهده دارد. سازوکار توجه از شواهد محلی برای افزایش نمایش ویژگی ها و استدلال در تولید توصیفات تصویری بهره برده و سازوکار توجه بر توجه می تواند روابط اشیا داخل تصاویر را به خوبی درک کند. این روش پیشنهادی توانسته است بر روی تصاویر مجموعه های داده Flickr8k و MSCOCO توصیف های بهتری را نسبت به روش های موفق موجود ارائه دهد. همچنین بر اساس معیارهای ارزیابی METEOR، ROUGE عملکرد توصیف تصویر را بهبود داده است.

### کلید واژه ها

توصیف تصویر، رمزگذار-رمزگشا سازوکار توجه، سازوکار توجه بر توجه، پردازش زبان ها طبیعی

### ۱- مقدمه

آموزش داده می شوند، می توانند برای ارائه اطلاعات مفید به پزشکان و کمک به روش های تشخیص استفاده شوند [۱]. سایر کاربردهای توصیف تصویر، شامل جستجو و بازیابی جملات تصویری و همچنین آوردن هوش بصری برای ربات ها است. در تولید خودکار توصیف متنی برای تصاویر در مقایسه با طبقه بندی تصویر، تشخیص اشیا و بسیاری دیگر از وظایف بینایی کامپیوتر، با چالش های بیشتری مواجه است. این مسئله به این دلیل است که برای تولید توصیفی که از لحاظ محتوایی غنی و از لحاظ نحوی صحیح باشد، نیاز به ادراک کلی و همه جانبه از موجودات و اشیا داخل تصویر و روابط بین آنهاست [۲]. همانطور که ذکر شد، آموزش یک مدل کارآمد توصیف تصاویر مستلزم درک جامعی از موجودیت های اساسی در تصویر و همچنین روابط آنها بوده و بسیار چالش برانگیز است. در این زمینه موفقیت های زیادی در

توصیف تصویر، پیوندی بین پردازش زبان طبیعی و بینایی ماشین ایجاد می کند. سیستم های توصیف تصویر را می توان برای کارهای مختلف استفاده کرد؛ به عنوان مثال، این سیستم ها می توانند افراد نابینا را قادر به دریافت اطلاعات بصری در مورد محیط اطراف خود کنند. هنگامی که این سیستم ها بر روی تصاویر پزشکی

این مقاله در آبان ماه ۱۴۰۱ دریافت، در بهمن ماه بازنگری و سپس پذیرفته شد.

آزمایشگاه هوش و بینایی ربات، گروه مهندسی کامپیوتر، دانشگاه بوعلی سینا  
رایانامه: [khotanlou@basu.ac.ir](mailto:khotanlou@basu.ac.ir)

نویسنده مسئول : حسن ختن لو

بیشتر روش‌های توصیف تصویر تلاش می‌کنند که مقادیر معیارهای BLEU رو افزایش دهند که شباهت بین دو جمله را حساب می‌کند. این معیار در واقع به‌عنوان میانگین هندسی حاصل ضرب دقت n-gram در جریمه مختصر بودن برای جملات کوتاه تعریف شده است. این معیار چون به هم معنی‌ها توجه نمی‌کند، معیار خوبی برای سنجش نیست. معیار METEOR یکی دیگر از معیارهای مورد استفاده است که این معیار به‌عنوان میانگین هارمونیک دقت، recall و تطابق unigram بین جملات می‌شود. علاوه بر این، در این معیار از مترادف‌ها و ریشه کلمات نیز استفاده می‌شود. این معیار ارزیابی چند نقص BLEU را مانند ارزیابی recall و عدم تطابق صریح کلمات، برطرف می‌کند. به همین دلیل معیار METOR نسبت به معیارهای ارزیابی دیگر، معیار بهتری برای سنجش مدل پیشنهادی است.

حال با بیان ویژگی‌های مسئله توصیف متنی تصاویر، چالشها و نمونه‌هایی از تلاشهای انجام شده در این زمینه، به معرفی روش پیشنهادی خود برای حل این مسئله می‌پردازیم. دستاوردهای این تحقیق به طور کلی به صورت زیر خلاصه می‌شود:

- ارائه مدلی بر مبنای معماری رمزگذار - رمزگشا با تلفیق سازوکارهای توجه بر توجه در جهت بهبود عملکرد مدل.
  - استفاده از سازوکار  $AoA^3$  و Multi-Head Attention برای نشان دادن موثرتر ویژگی‌ها و درک بهتر روابط میان اشیا.
  - مقایسه عملکرد با روش‌های موجود بر روی مجموعه‌های داده Flickr8k، MSCOCO با معیارهای ارزیابی BLEU(1-4)، METEOR، ROUGE، CIDER.
- در این بخش برای ایجاد یک فضای روشن از مطالعات انجام شده در جهت حل مسئله توصیف تصاویر مرتبط با روش پیشنهادی، پژوهش‌ها در دو بخش روشهای مبتنی بر معماری رمزگذار-رمزگشا و مبتنی بر سازوکار توجه، مرور و بررسی می‌شوند.

## ۲- کارهای مرتبط

### ۲-۱- روش‌های مبتنی بر معماری رمزگذار-رمزگشا

از آنجایی که استفاده از معماری رمزگذار - رمزگشا در ترجمه ماشینی باعث بهبود شده است، به همین دلیل روش‌های توصیف تصویر نیز از این روش استفاده می‌کنند. معماری رمزگذار - رمزگشا به طور معمول دو مرحله دارد: در مرحله اول، مدل از شبکه‌های کانولوشنی (CNN) به منظور استخراج ویژگی‌ها، تشخیص اشیا و ارتباط میان آن‌ها استفاده می‌کند. در مرحله دوم، از یک مدل زبانی استفاده می‌شود که خروجی مرحله اول را به شکل کلمات تولید می‌کند [۲]. ویلیناس و همکاران [۳] مدلی به نام

سال‌های اخیر به دست آمده، اما هنوز مشکلات زیادی وجود دارد. یکی از چالش‌های مهم این است که بیشتر تحقیقات بر روی شبکه‌های بازتابی تولید جمله تمرکز می‌کنند. آن‌ها تأثیر ویژگی‌های مشتق شده از تصویر را نادیده گرفته‌اند [۳]. چالش دیگر در مورد معیارهای ارزیابی پیشنهادی محققان است، که نمی‌توانند عملکرد مدل‌ها را در مقایسه با انسان به دقت اندازه‌گیری کنند [۴].

توصیه‌برچسب تصویر، با هدف تخصیص مجموعه‌ای برچسب‌های مرتبط برای تصاویر از جمله کارهایی بود که در این زمینه انجام شده است. مدل‌های مبتنی بر چسب بر این تمرکز می‌کند که چه قدر برچسب‌های پیشنهادی محتوای تصویر را توصیف می‌کند. از جمله روش‌هایی که در این زمینه انجام شده است، روش‌هایی مبتنی بر خوشه بندی است [۵، ۶]. اما هدف از این مقاله تولید توصیفات است که تصویر را شرح دهد.

تولید الگو و پرکردن اسلات [۷-۹] و بازتابی عنوان [۱۰-۱۳] از اولین تکنیک‌هایی هستند که برای تولید توصیف تصویر به صورت خودکار استفاده شده است. سپس، نتایج بهتری با استفاده از شبکه‌های عصبی عمیق در مقایسه با روش‌های مبتنی بر الگو و بازتابی حاصل شده است. روش‌های مبتنی بر شبکه‌های عصبی عمیق [۱۴-۱۶] به دلیل آموزش با تصاویر متعدد، بدون نیاز به طراحی ویژگی‌ها توسط انسان، دقت بالایی برای تصاویر ناشناخته دارند. یکی از رایج‌ترین معماری‌های عمیق مورد استفاده در تولید توصیف تصویر به صورت خودکار، معماری رمزگذار - رمزگشا<sup>۱</sup> است. مدل رمزگذار - رمزگشا معمولاً برای ارائه توصیف متنی از تصویر استفاده می‌شود. اطلاعات در سطح پیکسل به رمزگشای شبکه عصبی کانولوشنی<sup>۲</sup> (CNN) داده می‌شود تا با مترام کردن آن را رمزگذاری کند. سپس، رمزگشا برای ترجمه این اطلاعات به زبان‌های طبیعی استفاده می‌شود [۲]. مدل‌های دیگر نیز بر پایه معماری رمزگذار - رمزگشا. شبکه‌های بازگشتی برای استخراج ویژگی‌های تصویر و تولید عنوان پیشنهاد شدند [۳، ۱۷]. با توجه به شکاف زیاد بین بینایی و روش‌های پردازش زبان، اکثر روش‌ها مشکل عدم تطبیق معنایی کامل بین تصاویر و توصیف‌های تولید شده را دارند. برای حل این مشکل، روش‌های معنایی معرفی شدند. روش‌هایی مانند تطبیق معنایی هدایت شده با شباهت چند سطحی که می‌توانست شباهت‌های معنایی محلی و سراسری را برای یادگیری همبستگی پنهان بین تصاویر و توصیفات تولید شده ترکیب کند [۱۸]، از روش‌های برخورد این چالش بود. اخیراً، سازوکار توجه در نقشه ویژگی CNN مورد بررسی قرار گرفته است که می‌تواند یک نقشه فضایی ایجاد کند تا اهمیت یا برجستگی خاصی به مناطق مهم تصویری مرتبط با هر کلمه تولید شده بدهد [۱۹]. بنابراین، مدل‌ها می‌توانند روابط بین ویژگی‌های تصویر خاص و کلمات مربوطه در توصیف‌ها را بیاموزند [۱۷].

<sup>۱</sup>Encoder-Decoder

<sup>۲</sup>Convolution neural network(CNN)

<sup>۳</sup>Attention on attention

## ۲-۲- روش‌های مبتنی بر سازوکار توجه

در معماری رمزگذار-رمزگشا بدون استفاده از سازوکار توجه<sup>۴</sup>، یک CNN به عنوان رمزگذار برای استخراج ویژگی‌های بصری از تصویر و RNN<sup>۵</sup> برای تبدیل ویژگی‌های استخراج شده به کلمات استفاده شد. این روش نمی‌تواند تصویر را در طول زمان تجزیه و تحلیل کند و همچنین نمی‌تواند جنبه‌های فضایی تصویر را که مربوط به قسمت‌های مهم تصویر است را در نظر بگیرد. در روش رمزگذار - رمزگشا، توصیف با استفاده از کل صحنه تولید می‌شود. سازوکار توجه می‌تواند به صورت پویا روی عناصر مختلف یا مناطق تصویر متمرکز شود. خو و همکاران [۲۵] در سال ۲۰۱۵، اولین روش مبتنی بر سازوکار توجه برای تولید توصیف تصویر را معرفی کردند. تفاوت اصلی روش‌های مبتنی بر توجه با روش‌های دیگر این است که می‌تواند روی بخش‌های برجسته داده‌ها تمرکز کند و همزمان کلمات مرتبط را تولید کند. این روش از دو تکنیک مختلف استفاده می‌کند: توجه سخت تصادفی<sup>۶</sup> و توجه نرم قطعی. اکثر رویکردهای مبتنی بر CNN از لایه‌ی بالایی ConvNet برای استخراج اطلاعات برجسته شی از تصویر استفاده می‌کنند. پادرسولی و همکاران [۲۶] سازوکارهای توجه مبتنی بر منطقه را اعمال کردند. در این روش روابط بین کلمات تولید شده و نواحی تصویر مدلسازی شد. دنگ و همکاران [۲۷] در سال ۲۰۲۰ روشی را پیشنهاد کردند که از مدل DenseNet برای استخراج ویژگی‌های کلی تصویر در مرحله رمزگذاری استفاده می‌کند. و در همان زمان، سازوکار توجه از یک دروازه نگهبان<sup>۷</sup> استفاده می‌کند تا تصمیم بگیرد که از کدام اطلاعات و ویژگی تصویر برای تولید کلمه استفاده شود. آنها از LSTM برای بخش رمزگشایی مدل استفاده می‌کنند. این روش توانست توصیفات را تا حدودی بهبود بخشد. کائو و همکاران [۲۸] در سال ۲۰۲۰ مدل IGGAN را برای یادگیری بدون نظارت پیشنهاد کردند، که نمایش خصوصیات چند مقیاسی و روابط شی را ترکیب می‌کند. برای داشتن یک نمایش قدرتمند، تصاویر توسط ResNet با یک ماژول چند مقیاسی جدید و کانال RMCNET رمزگذاری می‌شود. روابط شی نیز توسط ماژول MAN استخراج می‌شود. شبکه IGGAN ترکیبی از شبکه‌های مولد و سازوکار توجه است. دپنگ و وانگ و همکاران [۲۹] روشی مبتنی بر معماری رمزگذار - رمزگشا را طراحی کردند که شامل یک رابط هدایت کننده متن<sup>۸</sup> بود تا نمایش بصری را یاد بگیرد که با شناخت بصری انسان سازگار تر باشد. علاوه بر این بخش رمزگشا را به دو ماژول تقسیم می‌کند: یک تولید کننده برای تولید جمله اولیه و یک پالایش کننده برای اصلاح

NIC<sup>۱</sup> پیشنهاد کردند. این روش از CNN برای نمایش ویژگی‌های تصویر و از LSTM<sup>۲</sup> برای تولید توصیف استفاده می‌کند. در این روش CNN از روش جدیدی برای نرمال سازی دسته‌ای بهره می‌برد. و خروجی آخرین لایه پنهان CNN به عنوان ورودی به رمزگذار داده می‌شود. این LSTM می‌تواند اشیایی را که قبلاً با استفاده از متن توضیح داده شده‌اند را ردیابی کند. عملکرد NIC بر اساس حداکثر برآوردهای آموزشی است. از آنجایی که اطلاعات تصویر فقط در ابتدای فرآیند تغذیه می‌شود، ممکن است با مشکل ناپدید شدن گرادیان مواجه شوند. به این ترتیب نقش کلمات تولید شده در ابتدای جملات در تولید کلمات بعدی به تدریج ضعیف و ضعیف تر می‌شود. بنابراین، LSTM هنوز با چالش‌هایی از جمله ایجاد جملات طولانی [۲۰-۲۲] مواجه است، به همین دلیل جیا و همکاران [۴] توسعه‌ای از LSTM به نام LSTMg را پیشنهاد کردند که می‌تواند جملات طولانی تولید کند. در این معماری، اطلاعات معنایی کلی با استفاده از دروازه و حالت سلول به LSTM اضافه می‌شود. وانگ و همکاران [۲۳] یک روش عمیق مبتنی بر LSTM دو طرفه پیشنهاد کردند. که این روش قادر است توصیف‌های غنی تولید کند. روش پیشنهاد شده از یک CNN و دو شبکه مجزا LSTM تشکیل شده است.

برای بهبود روش‌های مبتنی بر چارچوب رمزگذار-رمزگشا، محققان روش‌هایی مبتنی بر فضای چندوجهی پیشنهاد کردند. معماری فضای چند وجهی<sup>۳</sup> شامل سه بخش رمزگذار زبان و بینایی، فضای چند وجهی و رمزگشای زبان است. اولین کار در این زمینه توسط کیروس [۱۶] در سال ۲۰۱۴ انجام شده است. این روش از CNN برای استخراج ویژگی‌های تصویر و همچنین از فضای چند وجهی که تصویر و متن را نشان می‌دهد، استفاده می‌کند. چن در سال ۲۰۱۵ [۲۴] مدل توصیف تصویر را بر اساس فضای چند وجهی پیشنهاد کرد. این روش می‌توانست توصیف‌های جدیدی برای تصویر ایجاد کند، و ویژگی‌های بصری را از توصیف‌های داده شده بازیابی و همچنین نمایش بصری تصاویر را از کلمات تولید شده به صورت پویا به روزرسانی کند. این مدل برای کار با داده‌های زیاد محدودیت‌های فراوانی دارد و در کار با حافظه بلند مدت ناکارآمد است. به دلیل وجود چالش‌های توضیح داده شده، در مرحله بعد، سازوکار توجه در این چارچوب به کار گرفته شد که تولید توصیف برای تصویر را بهبود بخشد.

<sup>۴</sup>Attention Mechanism

<sup>۵</sup>شبکه‌های بازگشتی

<sup>۶</sup>Random hard attention

<sup>۷</sup>Sentinel gate

<sup>۸</sup>Text-guided relation encoder(TGER)

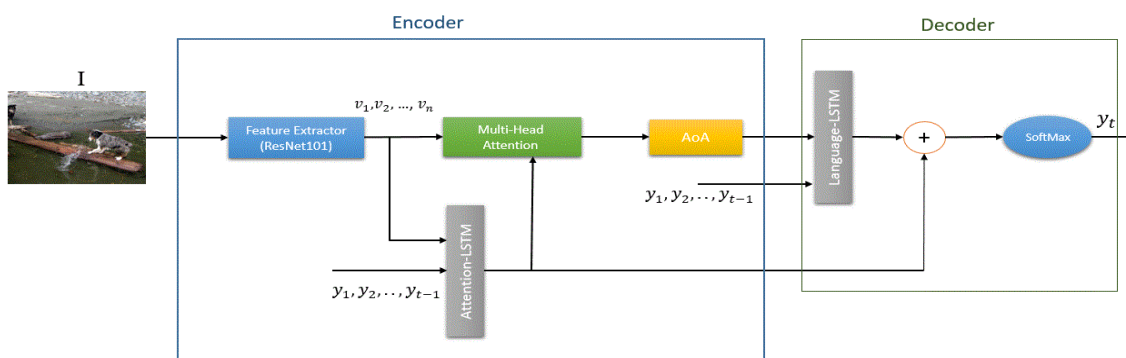
<sup>۱</sup>Generator Neural Image Caption

<sup>۲</sup>Long short-term memory

<sup>۳</sup>Multi Modal Space

را به صورت منطقی و روان تولید کند. پالایش کنند هاشاملیک  
ماژول رمزگذار توصیف، یک LSTM

جمله. ماژول تولید کننده متشکل از یک LSTM استاندارد، و یک  
ماژول Gate on Attention (GOA) است که سعی دارد جمله اصلی



شکل (۱): ساختار کلی روش پیشنهادی

نشان داده شوند. همچنین در نمایش جملات، طول توضیحات  
تولید شده را نشان می‌دهد. همانطور که در معادله (۱) مشاهده  
می‌کنید،  $I$  به عنوان ورودی تصویر به مدل داده می‌شود و  
توضیحات  $S$  در خروجی ایجاد می‌شود.

$$F: I \rightarrow S \quad (1)$$

روش پیشنهادی ما بر اساس پژوهش وی و همکاران [۳۱]،  
لان هانگ و همکاران [۳۰] و مطالعه پیشینمان [۲] است. ما در  
روش پیشین خود نیز از معماری رمزگذار و رمزگشا استفاده کرده  
ایم که شامل چندین ماژول استخراج ویژگی، Attention-LSTM،  
Attention-Layer، Language-LSTM است. در این معماری  
لغات از پیش تولید شده و ویژگی‌های استخراج شده به لایه  
Attention-LSTM داده می‌شود که این لایه شامل دو لایه LSTM  
است و سپس به یک لایه ساده توجه داده می‌شود، و سپس  
خروجی مدل به لایه Language-LSTM می‌رود که شامل یک لایه  
LSTM است و در نهایت کلمه تولید می‌شود.

معماری مطالعه پیشین در روش پیشنهادی تغییر کرده است و  
ضمن کاهش پیچیدگی محاسباتی نسبت به کار وی و همکاران،  
نتایج نیز بهبود یافتند. ما در روش پیشنهادی خود بر خلاف مطالعه  
پایه [۲۸] از شبکه‌های مولد استفاده نکرده‌ایم. ساختار Multi-  
Head Attention، Attention-LSTM نسبت به کار آنان تغییر کرده  
است، اما ایده ماژول AoA از مطالعه پایه گرفته شده و به مدل  
پیشنهادی اضافه شده است. در واقع ما در این روش، مدل قبلی  
پیشنهادی خود را با عوض کردن ساختار لایه سازوکار توجه و با  
افزودن لایه AoA بهبود داده‌ایم، ما از روش تخمین حداکثر  
احتمال<sup>۲</sup> (MLE) برای آموزش مدل استفاده کرده‌ایم. در این روش  
در زمان  $T$  یک توصیف واقعی تا کلمه  $T-I$  به مدل داده می‌شود تا  
احتمال تولید  $y_t$  به عنوان عنصر واقعی جمله مدل به حداکثر برسد.  
برای بهینه سازی مدل از تابع تلفات cross-entropy استفاده  
می‌کنیم که در معادله ۲ نمایش داده شده است.

مبتنی بر توجه و یک ماژول GOA است، که به طور مکرر جزئیات را  
در عنوان اصلی تغییر می‌دهد تا زیرنویس‌ها را با کمک متن راهنما  
غنی و دقیق تولید کند. لان هانگ و همکاران [۳۰] روشی را مبتنی  
بر روش توجه بر توجه (AoA) معرفی کردند که توانست توصیفات  
خوبی را تولید کند. در این روش AoA به قسمت رمزگذار و  
رمزگشای مدل افزوده شد. AoA یک گسترش کلی برای سازوکار  
توجه است و می‌تواند به هر کدام از آن‌ها اعمال شود. در قسمت  
رمزگذار، AoA به مدل سازی بهتر روابط بین اشیا مختلف در  
تصویر کمک می‌کند. در رمزگشا نیز AoA نتایج توجه غیر مرتبط  
را فیلتر می‌کند و فقط موارد مفید را نگه می‌دارد، همچنین ما  
روش [۲] را پیشنهاد کرده بودیم که از یک لایه توجه محلی<sup>۱</sup> در  
چارچوب رمزگذار - رمزگشا استفاده شد، که این روش باعث شد  
توصیف‌های خوبی را مدل پیشنهادی تولید کند و معیارهای  
ارزیابی را بهبود دهد و پژوهش حاضر توسعه یافته این روش می  
باشد.

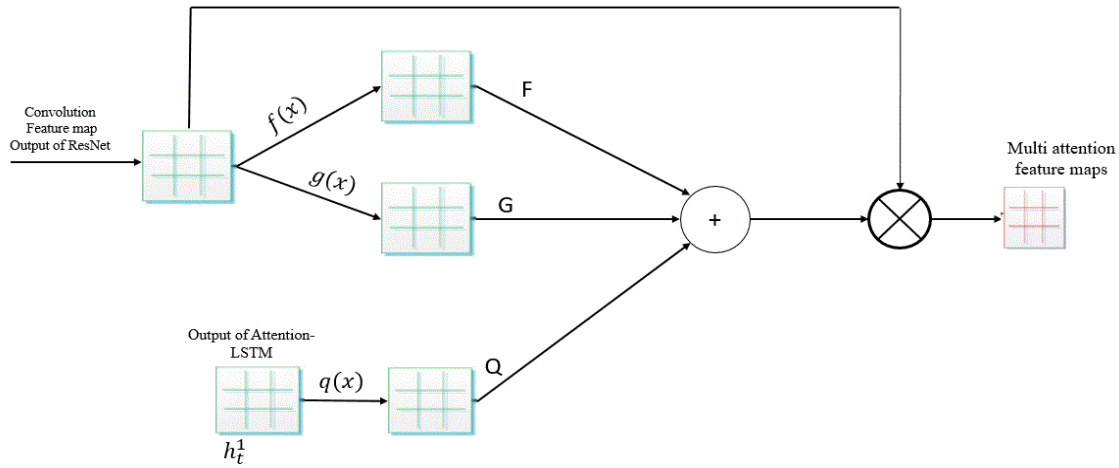
### ۳- روش پیشنهادی

در این قسمت ما روش پیشنهادی خود را معرفی کرده و جزئیات  
آن را بررسی می‌کنیم. ساختار کلی روش پیشنهادی در شکل (۱)  
مشاهده می‌شود. این ساختار از چارچوب رمزگذار - رمزگشا  
پیروی می‌کند و از چندین بخش تشکیل شده است. قسمت  
رمزگذار تصویر شامل چند ماژول استخراج کننده ویژگی  
(ResNet101، Multi-Head Attention، Language-LSTM،  
Attention-On-Attention است. قسمت رمزگشای تصویر شامل  
ماژول Language-LSTM است.

همانطور که در شکل یک نشان داده شده، تصویر  $I$  به عنوان  
ورودی به مدل داده می‌شود. مدل توصیف تصویر باید جملاتی را  
تولید کند که از لحاظ معنایی و محتوایی صحیح باشد و می‌تواند  
با نمایش  $S = y_1, y_2, \dots, y_L$  که در آن  $y_i$  کلمه  $i$ ام جمله  $S$  است،

<sup>۲</sup>Maximum likelihood estimation

<sup>۱</sup>Local attention mechanism



شکل (۲): ساختار توجه چند وجهی پیشنهادی

ادامه می‌دهد. همانطور که در شکل (۱) نشان داده شده است، پس از ادغام ویژگی‌های سراسری  $V$ ، با بردار کلمه از پیش تولید شده  $y_t$   $W_e$ ، به این لایه داده می‌شود و در ادامه چنان که در معادله (۴) مشاهده می‌کنید، خروجی به لایه LSTM داده می‌شود:

$$h_t^1 = LSTM(LSTM([v \cdot W_e y_{t-1}])) \quad (4)$$

که در معادله (۴)  $V = \{v_1, v_2, \dots, v_n\}$  نشان دهنده ویژگی‌های کلی استخراج شده از تصویر است که در قسمت قبل استخراج شده است.  $W$  بردار کلمات تعبیه شده از پیش تولید شده است.

### ۳-۱-۳ Multi-Head Attention

این ماژول شامل یک سازوکار چند توجهی با استفاده از شواهد محلی و غیر محلی برای بازنمایی و استدلال موثرتر ویژگی در توصیف تصویر پیشنهاد شده است. همچنین در این معماری، توجه یک لایه متراکم است که ویژگی‌های سطح پایین‌تر یک منطقه محلی را دستکاری می‌کند تا یک نقشه فضایی ایجاد کند که مناطق تصویر را بر اساس اطلاعات بازخورد ارائه شده توسط Attention-LSTM برجسته می‌کند. همانطور که در شکل (۱) مشاهده می‌کنید، خروجی لایه Feature-Extractor ( $V$ ) و خروجی لایه Attention-LSTM به عنوان ورودی به این لایه داده می‌شود. در شکل (۲) جزئیات این لایه را به صورت دقیق نمایش داده‌ایم. در ابتدا ویژگی‌های استخراج شده  $V$  به دو فضای  $F$ ،  $G$  منتقل می‌شوند. سپس خروجی لایه ( $Q$ ) Attention-LSTM به فضای جدید  $Q$  منتقل می‌شود. نمایش ریاضی این مراحل در معادلات ۵ تا ۷ نشان داده شده است.

$$f(v) = W_f V \quad (5)$$

$$g(v) = W_g V \quad (6)$$

$$q(h_t^1) = W_q h_t^1 \quad (7)$$

و سپس جمع  $\tanh$  را به صورت معادله ۸ محاسبه می‌کنیم:

$$L_{CE}(\theta) = -\sum_{t=1}^T \log(P_{\theta}(y_t^* | v, y_1^*, \dots, y_{t-1}^*)) \quad (2)$$

پارامتر  $\theta$  وزن‌های قابل یادگیری مدل را نشان می‌دهد. این مدل توزیعی از وجود هر کلمه را به شرط دنباله‌ای از پیش موجود از کلمات و ویژگی‌های بصری ایجاد می‌کند؛ زیرا هیچ توصیف واقعی از فرم موجود وجود ندارد.

برای شرح عملکرد مدل، مطابق شکل (۱)، ایجاد یک عنوان برای تصویر که از نظر نحوی و معنایی صحیح باشد را مطابق مراحل زیر دنبال می‌کنیم:

### ۳-۱-۳-۱ رمزگذار

قسمت‌های رمزگذار تصویر در مدل پیشنهادی شامل: استخراج کننده ویژگی‌های تصویر، Attention-LSTM، Multi-Head Attention، Attention On Attention است، که در ادامه توضیح خواهیم داد.

### ۳-۱-۱-۱ استخراج کننده ویژگی

ویژگی‌های کلی تصویر، یک تصویر را به شکل جامع توصیف می‌کند. ویژگی‌های کلی تصویر محبوب هستند، زیرا نمایش بسیار فشرده‌ای از تصاویر ایجاد می‌کنند این مدل در مجموعه داده ImageNet از قبل آموزش داده شده است. ویژگی‌های تصویر کلی  $V = \{v_1, v_2, v_3, \dots, v_k\}$  جزئیات بیشتری را نسبت به تصویر ورودی پوشش می‌دهد.  $v_i$  یک بردار  $d$ -بعدی است که ویژگی‌های یک ناحیه از تصویر را نشان می‌دهد و  $k$  تعداد مناطق تصویر است.

$$V = CNN_{CONV}(I) \quad (3)$$

در فرمول (۳)  $I$  تصویر ورودی است و  $CNN_{CONV}$  آخرین لایه کانولوشن ResNet101 است و در نهایت  $V$  ویژگی‌های کلی تصویر است که استخراج می‌شود.

### ۳-۱-۲ Attention-LSTM

این ماژول به مفاهیم مرتبط به زبان و اطلاعات متنی از پیش تولید شده توجه می‌کند. سپس به توزیع توجه در تمام قسمت‌های تصویر

## ۴- ارزیابی

### ۴-۱- مجموعه داده

Flickr8k [۱۳] مجموعه داده‌ی گسترده‌ای برای توصیف تصاویر است. Flickr8k یک پایگاه داده از ۸۰۰۰ تصویر است که هر کدام دارای پنج عنوان به زبان طبیعی است که بر اساس سرویس مجموعه



Caption1: A black dog is running after a white dog in the snow.  
Caption2: Black dog chasing brown dog through snow.  
Caption3: Two dogs chase each other across the snowy ground.  
Caption4: Two dogs play together in the snow.  
Caption5: Two dogs running through a low lying body of water.

Caption1: A group of cyclists race uphill on foot with their bikes on their shoulders.  
Caption2: Bicyclists are carrying their bicycles up a hill.  
Caption3: Cyclists are carrying their bikes up the hill.  
Caption4: people carry their bikes up the hill.  
Caption5: Several bicyclists are carrying their bikes up a grassy hill.

شکل (۴): نمونه‌ای از مجموعه داده

منابع آمازون مکانیکال ترک تولید شده‌اند. تصاویر، بیشتر در مورد انسان و حیوانات بوده و هر عنوان یا توصیف بخشی از موجودات و رویدادهای برجسته تصویر مربوطه را ارائه می‌دهد. شکل (۴) نمونه‌ای از مجموعه داده Flickr8k را نشان می‌دهد. مجموعه داده [۳۲] نیز مجموعه داده جامعی با ۱۲۳۲۸۷ تصویر است که در آن برای هر تصویر پنج توصیف ارائه شده است.

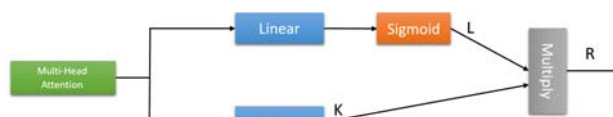
### ۴-۲- معیار ارزیابی

در حال حاضر معیارهای ارزیابی تولید توصیف تصویر همان معیارهای ارزیابی ترجمه ماشینی است. زیرا روشهای ارزیابی انسانی منسوخ شده و توصیف تصویر معیارهای ارزیابی جداگانه‌ای ندارد. مشکل معیارهای ارزیابی موجود این است که شباهت و عدم شباهت بین توصیف تولید شده و توصیف‌های مرجع را نشان می‌دهد و فقط به جنبه زبانی توجه می‌کنند، اما در حال حاضر تنها روش ارزیابی همین معیارهاست. معیارهای ارزیابی خودکار که به صورت گسترده استفاده می‌شود شامل: BLEU, METEOR, CIDER, ROUGE است. در این معیارها، نمرات بالاتر، جملات بهتر را نشان می‌دهد.

معیار BLEU [۳۳]، یکی از اولین معیارهای ارزیابی است که برای اندازه‌گیری شباهت بین دو جمله مورد استفاده قرار گرفته است. این معیار در ابتدا برای ترجمه ماشینی پیشنهاد شده و به عنوان میانگین هندسی، حاصل ضرب دقت n-gram در جریمه مختصر بودن برای جملات کوتاه تعریف شده است. نحوه محاسبه دقت n-gram در فرمول (۱۴) ذکر شده است.

$$P_n = \frac{\sum_{C \in \{Candidates\}} \sum_{n-grams \in C} Count_{clip}(n-gram)}{\sum_{C \in \{Candidates\}} \sum_{n-gram \in C} Count(n-gram)} \quad (15)$$

$$a_{i,t} = \tanh(f + g + q) \quad (8)$$



شکل (۳): ساختار پیشنهادی توجه بر توجه

$a_{i,t}$  وزن‌های استخراج شده از لایه Multi Head Attention است.

و در مرحله آخر در این لایه وزن بدست آمده  $a_{i,t}$  را ضرب در ویژگی‌های استخراج شده  $V$  از ResNet می‌کنیم.

$$c_t = \sum_{k=1}^k a_{k,t} v_k \quad (9)$$

### ۳-۱-۱- Attention on Attention

AoA یک گسترش کلی برای سازوکار توجه است. و می‌تواند به هرکدام از آن‌ها اعمال شود. در قسمت رمزگذار، AoA به مدل سازی بهتر روابط بین اشیا مختلف در تصویر کمک می‌کند. همانطور که در شکل (۳) مشاهده می‌کنید، خروجی لایه Multi-Head Attention با نام  $C$  به عنوان ورودی به این لایه داده می‌شود. که از لایه خطی عبور داده می‌شود، در ابتدا  $C$  از لایه خطی  $W_l$  عبور داده می‌شود و سپس sigmoid آن محاسبه می‌شود و همینطور  $C$  از لایه خطی  $W_k$  عبور داده می‌شود و در نهایت خروجی آن‌ها در هم ضرب می‌شود. معادلات ۱۰ تا ۱۲ توصیف مراحل گفته شده است:

$$l(c) = \text{Sigmoid}(W_l C) \quad (10)$$

$$k(c) = W_k C \quad (11)$$

$$R = l * k \quad (12)$$

### ۳-۲- رمزگشا

#### ۳-۲-۱- Language-LSTM

همانطور که در شکل (۱) مشاهده می‌کنید، ورودی این ماژول شامل خروجی تولید شده توسط ماژول AoA و کلمات از پیش تولید شده تعبیه شده می‌باشد، که به شرح معادله (۱۳) است:

$$h_t^2 = \text{LSTM}(l r_t, W_e y_{t-1}) \quad (13)$$

در نهایت، مجموع Attention-LSTM و Language-LSTM به یک تابع softmax که در معادله ۱۴ آمده، داده می‌شود تا توزیع احتمال کلمه را بدست آوریم:

$$p(y_t | y_{1:t-1}) = \text{softmax}(h_t^1 + h_t^2) \quad (14)$$

که طولانی‌ترین دنباله‌های مشترک بین جفت جمله را اندازه‌گیری می‌کند [۳۶].

#### ۴-۳- آماده سازی داده

به دلیل ناکافی بودن داده‌ها در مجموعه داده Flickr8k، ما به تقویت داده<sup>۱</sup> در این مجموعه نیاز داشتیم. از این رو با روشهای مختلف، ۲۰۰۰ داده به مجموعه داده Flickr8k اضافه کردیم. ما داده‌ها را به روش‌های مختلفی تقویت کردیم. از جمله: افزودن نویز گاوسی<sup>۲</sup> به تصاویر، چرخاندن تصاویر از چپ به راست<sup>۳</sup>، افزودن روشنایی<sup>۴</sup> به تصاویر و آینه کردن تصاویر.

#### ۴-۴- جزئیات پیاده سازی

برای پیاده سازی مدل پیشنهادی، از زبان برنامه نویسی پایتون نسخه ۳٫۹ در چارچوب کتابخانه Keras استفاده کردیم. در این تحقیق از روشی برای تفکیک داده‌های ورودی به داده‌های آموزشی و داده‌های آزمون استفاده شد. که ۸۰ درصد داده‌های ورودی برای آموزش مدل و ۲۰ درصد برای آزمون مدل تقسیم شدند. همچنین یک کامپیوتر با یک پردازنده ۶۷۰۰-۱۷ Core شرکت اینتل با ۲۴ گیگابایت حافظه رم و یک پردازنده گرافیکی ۳۰۷۰ NVIDIA GTX مورد استفاده قرار گرفت. در روش پیشنهادی، تعداد حالت‌های پنهان LSTM و تعبیه ۱۲۸ است. همچنین از مدل از پیش آموزش دیده شده ResNet-101 در شبکه تصویر استفاده کردیم. تعداد پارامترهای مدل ما برابر با داده‌ایم. مدل ما بر مجموعه داده Flickr8k در ۱۰۰ دوره آموزش دید و همچنین این آموزش حدود ۹۸ ساعت طول کشید. همچنین آموزش مدل بر مجموعه داده MSCOCO در ۱۱۰ دوره آموزش حدود ۱۸۰ ساعت طول کشید. برای محاسبه‌ی معیارهای ارزیابی از کتابخانه‌ی استاندارد pycocoevalcap استفاده کرده‌ایم.

#### ۴-۵- نتایج تجربی

ما در این روش سعی بر تولید توصیفاتی داشتیم که از لحاظ محتوایی غنی و از لحاظ نحوی صحیح باشد. همانطور که قبلاً ذکر شد، یکی از بزرگترین مشکلات موجود در توصیف خودکار تصویر عدم وجود معیار ارزیابی مناسب است. معیارهای ارزیابی موجود و همچنین روش‌های پیشین سعی بر این دارند که توصیفی تولید کنند که بسیار شبیه به جمله مرجع باشد. اما ما در این روش سعی داشتیم روشی را پیشنهاد کنیم، که طی آن جملاتی تولید شوند که ضمن اینکه دقیقاً جمله مرجع نباشند اما به خوبی تصویر

که در اینجا  $\text{Count}_{\text{clip}}$  حداکثر تعداد n-gram است که در جمله کاندید و مرجع هم زمان اتفاق می‌افتد،  $\text{Count}$  تعداد n-gram‌هایی است که در جمله کاندید است. برای جلوگیری از ترجمه‌های بسیار کوتاه جریمه مختصر بودن را نیز اضافه کرده‌اند.

$$BP = \begin{cases} 1 & \text{if } |c| > |r| \\ e^{-(1-|r|/|c|)} & \text{if } |c| \leq |r| \end{cases} \quad (16)$$

در معادله ۱۶،  $|c|$  طول جمله کاندید و  $|r|$  طول جمله مرجع است. و در نهایت BLEU به صورت فرمول (۱۷) محاسبه می‌شود که  $w_n$  در آن ضریب وزنی است.

$$BLEU = BP \cdot \exp\left(\sum_{n=1}^N w_n \log P_n\right) \quad (16)$$

METEOR [۳۴] معیار دیگری است که به عنوان میانگین هارمونیک دقت، recall و تطابق یک-گرام بین جملات تعریف می‌شود. همچنین در این معیار از مترادف و ریشه کلمات استفاده می‌کند. این متریک چندین نقص BLEU، مانند عدم تطابق کلمات صریح را برطرف می‌کند. آنها همچنین قادر به شناسایی شباهت معنایی بودند که توسط تطبیق مترادف مبتنی بر WordNet کنترل می‌شود. تعداد unigram نگاشت شده بین دو رشته  $m$  نام دارد و تعداد unigram در جمله کاندید  $t$  و همچنین تعداد unigram مرجع  $r$  نام‌گذاری شده است. دقت unigram برابر است با  $P=m/t$  و همچنین recall آن برابر است با  $R=m/t$ ، سپس میانگین هارمونیک  $P$  و  $R$  مطابق معادله ۱۷ محاسبه می‌شود.

$$F_{mean} = \frac{P \cdot R}{\alpha \cdot P(1-\alpha) + R} \quad (17)$$

همچنین از تطابق n-gram طولانی‌تر برای محاسبه استفاده می‌شود. و در نهایت معیار METEOR طبق فرمول (۱۸) محاسبه می‌شود همچنین جریمه  $P$  برای تراز استفاده می‌شود.

$$M = F_{mean}(1-p) \quad (18)$$

هدف معیار ارزیابی CIDER [۳۵] محاسبه تناظر بین جمله نامزد  $c_i$  و مجموعه توصیفات مرجع تصویر به صورت  $S_i = S_{i1}, \dots$  است. برای محاسبه این معیار، همه کلمات در جملات مرجع نامزد به ریشه نگاشت می‌شوند. سپس همزمانی رخداد n-gram در جملات مرجع و نامزد محاسبه می‌شود. در CIDER، مشابه tf-idf n-gram+idf که در تمام توضیحات تصویر مشترک هستند، سبک وزن هستند. در نهایت، شباهت کسینوس بین n-gram گرم جملات کاندید و مرجع محاسبه می‌شود. این معیار گاهی به جزئیات بی اهمیت یک جمله وزن بیشتری می‌بخشد و در نتیجه ارزیابی عنوان بی اثر می‌شود [۳۶].

$$CIDER(c_i, s_i) = \sum_{n=1}^N w_n CIDER_n(c_i, c_n) \quad (18)$$

که  $CIDER_n$  برای n-gram به طول n شباهت کسینوسی بین جمله مرجع و کانید است.

معیار ارزیابی ROUGE [۳۷] در ابتدا برای ارزیابی سیستم‌های خلاصه سازی پیشنهاد شده است. و این ارزیابی از طریق مقایسه n-gram‌های همپوشانی، دنباله‌های کلمه و جفت کلمات انجام می‌شود. در این کار ما از L-Rouge استفاده می‌کنیم

<sup>۱</sup>Data augmentation

<sup>۲</sup>Gaussian augmentation

<sup>۳</sup>Filip augmentation

<sup>۴</sup>Brightness augmentation

توصیف می‌کند. همانطور که مشاهده می‌کنید این جملات دقیق شبیه به جمله مرجع نیست، به همین دلیل باعث شده است معیارهای BLEU که هدف از آن محاسبه‌ی شباهت n-gram است، کاهش بیابد. به طور مثال همانطور که در تصویر مشاهده می‌کنید در جمله مرجع کلمه woods وجود ندارد، اما روش ما به دلیل استفاده از مکانیزم توجه توانسته است woods را در تصویر تشخیص دهد. و در تصویر دوم در شکل (۶)، کلمه is running به جای کلمه jumps تشخیص داده شده است، این باعث می‌شود معیار BLEU کاهش پیدا کند در حالیکه جملات تولید شده از لحاظ معنایی و نحوی صحیح است.

بنابراین همانطور که مشاهده می‌کنید در جداول مقایسه روش پیشنهادی با روش‌های پیشین در معیارهای BLUE افزایش نداشته‌ایم و در معیار METEOR و ROUGE بهبود خوبی داشته‌ایم.

در جدول (۲)، مقایسه بین مطالعات موفق انجام شده و روش ما روی مجموعه داده Flickr8k نشان داده شده است. مدل پیشنهادی ما توانسته است مقادیر خوبی را در معیارهای ارزیابی بدست آورد. و همچنین در معیار METEOR که قبلاً اهمیت آن شرح داده شد، پیشرفت خوبی اتفاق افتاده است. همانطور که گفته شد، یکی از بزرگترین مشکلات در زمینه تولید توصیف تصاویر، خوب نبودن معیارهای ارزیابی است، بنابراین علاوه بر مشاهده نتایج معیارهای ارزیابی با تجزیه و تحلیل توضیحات تولید شده، می‌توان فهمید که توضیحات خوبی توسط مدل پیشنهادی ما تولید شده است. در جدول (۳) نیز مقایسه روش پیشنهادی با روش‌های پیشین بر مجموعه داده MSCOCO انجام شده است، همانطور که در جدول مشاهده می‌کنید در این مجموعه داده نیز مدل پیشنهادی توانسته است مقادیر خوبی را در معیارهای ارزیابی بدست آورد و همچنین معیار ROUGE نیز پیشرفت خوبی داشته است.

شکل (۷) نمونه‌ای از درک صحیح تصویر و اجزای آن توسط مدل ما است. در شکل (۸) مقایسه توصیف تولید شده توسط روش پیشنهادی و توسط یکی از روش‌های پیشین انجام شده است. هر دو مدل توانسته‌اند جملات خوبی را تولید کنند، اما مدل پیشنهادی MAGAN سعی بر این داشته که توصیفی تولید کند که بیشتر شبیه به توصیف مرجع باشد، اما روش ما توصیفی تولید کرده است که هم از لحاظ محتوایی صحیح هست و هم دقیقاً جمله مرجع نیست. به این ترتیب در روش ما ضمن تولید خروجی مناسب برای این تصویر از دایره واژگان گسترده‌تری استفاده شده که مطلوبیت و کیفیت مناسبی را از دید ناظر انسانی به خروجی می‌دهد.

را بیان و توصیف کند و از لحاظ نحوی صحیح باشد. به همین دلیل معیارهای ارزیابی موجود معیار مناسبی برای سنجش کار نبود. به دلیل اینکه معیار مناسب دیگری وجود نداشت معیارهای METEOR و ROUGE که نسبتاً معیارهای مناسب تری برای سنجش هستند برای ارزیابی انتخاب شدند. به خصوص معیار METEOR که از ریشه کلمات و هم معنی آن برای مقایسه استفاده می‌کند معیار مناسبی برای سنجش عملکرد مدل ما بود. در شکل (۵) نمونه‌ای از توصیف تولید شده بر مجموعه داده Flickr8k را مشاهده می‌کنید. همانطور که مشاهده می‌شود، این توصیف توانسته است، به خوبی تصویر را توصیف کند و همچنین از لحاظ نحوی صحیح است.



A man in a wetsuit rides his bicycle on a beach

شکل (۵) : نمونه‌ای از توصیف تولید شده توسط مدل پیشنهادی

#### ۴-۶- تجزیه و تحلیل کمی

در جدول (۱) مقایسه کار انجام شده، با روش قبلی پیشنهادی که روش پایه ما هست را نشان می‌دهد. در روش پیشین خود نیز از معماری رمزگذار و رمزگشا استفاده کرده ایم که شامل چندین ماژول استخراج کننده ویژگی، Attention-LSTM، Attention-Layer، Language-LSTM است. در روش پیشین [۲] نیز مدل پیشنهادی قبلی ما توصیفاتی تولید می‌کرد که از لحاظ نحوی و معنایی صحیح بود و توانسته بود در معیارهای ارزیابی مقادیر خوبی کسب کند، ما با بهبود لایه‌ی توجه و افزودن لایه‌ی توجه بر توجه به این مدل، همانطور که در جدول (۱) مشاهده می‌کنید، توانستیم معیارهای ارزیابی را بهبود دهیم.

در شکل (۶) نیز مقایسه توصیفات تولید شده توسط دو مدل را مشاهده می‌کنید. همانطور که می‌بینید، مدل جدید پیشنهادی ما توصیفات بهتری را تولید کرده است. اما با وجود اینکه توصیفات خیلی بهتر شده‌اند از دید معیارهای ارزیابی استفاده شده، تغییر چشم گیری نداده است. این مسئله تایید دیگری بر عدم ارزیابی کامل و دقیق توسط معیارهای موجود است. و همانطور که مشاهده می‌کنید، هر سه توصیفی که در شکل (۶) تولید شده است، از لحاظ معنایی و نحوی صحیح است، و به خوبی تصویر را



جدول (۱) : مقایسہ روش پیشنهادی جدید با روش پیشنهادی قبلی

| روش                    | Flickr8k |      |      |      |        |         |       |
|------------------------|----------|------|------|------|--------|---------|-------|
|                        | B-1      | B-2  | B-3  | B-4  | METEOR | ROUGE-L | CIDER |
| روش پیشین پیشنهادی [۲] | ۶۰,۴     | ۴۱,۳ | ۲۷,۷ | ۱۸,۷ | ۲۴,۴   | ۵۱,۶    | ۴۱,۵  |
| روش پیشنهادی           | ۶۱,۶     | ۴۲,۳ | ۲۷,۸ | ۱۸,۹ | ۲۵,۴   | ۵۲,۶    | ۴۲,۵  |



New-Method: a man is riding a bike on a trail through the woods.  
 Previous-Method: a person is riding a bike on a dirt bike.  
 Ground-Truth: a person riding a bike doing a very high jump.  
 Ground-Truth: a person is in the air with a trick bike



New-Method: a black dog is running through a field of snow.  
 Previous-Method: a black dog is running through the snow.  
 Ground-Truth: a black dog jumps in the snow.  
 Ground-Truth: a black dog bounds through a path in the snow.



New-Method: a basketball player dribbles the ball.  
 Previous-Method: a basketball player in a white uniform is playing with a ball in the background.  
 Ground-Truth: a basketball player dribbles the ball.  
 Ground-Truth: a basketball player in the white strip is dribbling the ball on the court.

شکل (۶) : مقایسہ توصیفات تولید شدہ

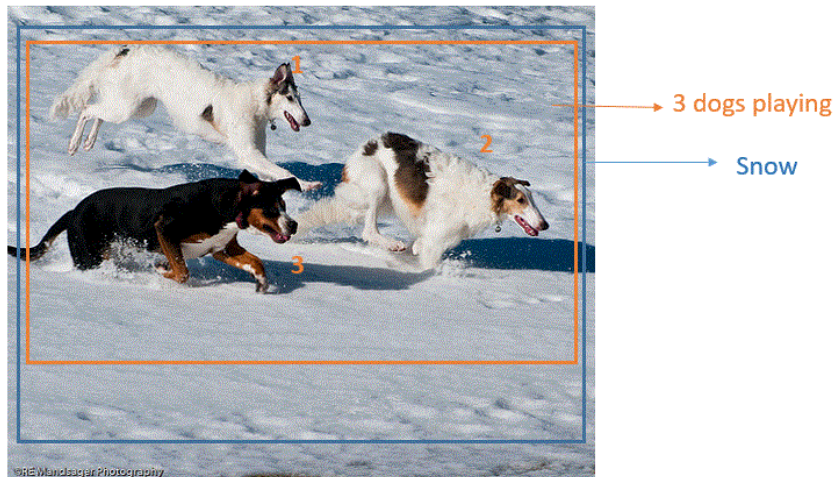
جدول (۲) : مقایسہ روش پیشنهادی جدید با روش های قبلی بر روی مجموعه داده Flickr8k

| روش                    | Flickr8k |      |      |      |        |         |       |
|------------------------|----------|------|------|------|--------|---------|-------|
|                        | B-1      | B-2  | B-3  | B-4  | METEOR | ROUGE-L | CIDER |
| خو و همکاران [۲۵]      | ۶۷       | ۴۵,۷ | ۳۱,۴ | ۲۱,۳ | ۲۰,۳   | -       | -     |
| فو و همکاران [۳۸]      | ۶۳,۹     | ۴۵,۹ | ۳۱,۹ | ۲۱,۷ | ۲۰,۴   | ۴۷      | ۵۳,۸  |
| لی و همکاران [۳۹]      | ۵۷,۲     | ۳۷,۹ | ۲۳,۹ | ۱۴,۸ | ۱۶,۶   | ۴۱,۹    | ۳۶,۲  |
| زو و همکاران [۴۰]      | ۶۳,۲     | ۴۴,۸ | ۳۱,۳ | ۲۱,۵ | ۲۰,۴   | -       | -     |
| ونگ و همکاران [۲۳]     | ۶۸,۹     | ۴۷,۴ | ۳۳,۴ | ۲۳,۸ | ۲۱,۴   | ۴۸,۸    | ۵۹,۴  |
| چنگ و همکاران [۴۱]     | ۶۳,۹     | ۴۶,۲ | ۳۲,۵ | ۲۲,۲ | -      | -       | -     |
| دینگ و همکاران [۴۲]    | ۶۷,۲     | ۴۵,۱ | ۳۰,۵ | ۲۱,۵ | -      | -       | -     |
| شیائو و همکاران [۴۳]   | ۶۴,۸     | ۴۶,۷ | ۳۲,۳ | ۲۱,۸ | ۲۰,۵   | -       | -     |
| روش پیشین پیشنهادی [۲] | ۶۰,۴     | ۴۱,۳ | ۲۷,۷ | ۱۸,۷ | ۲۴,۴   | ۵۱,۶    | ۴۱,۵  |
| روش ما                 | ۶۱,۶     | ۴۲,۳ | ۲۷,۸ | ۱۸,۹ | ۲۵,۴   | ۵۲,۶    | ۴۲,۵  |

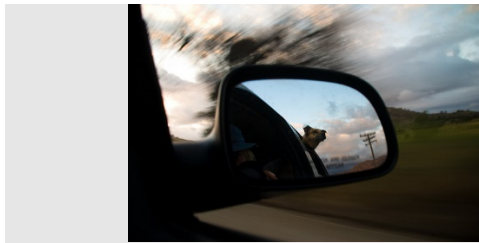
جدول (۳): مقایسه روش پیشنهادی جدید با روش‌های قبلی بر روی مجموعه داده MSCOCO

| روش                  | MSCOCO |      |      |      |             |             |       |
|----------------------|--------|------|------|------|-------------|-------------|-------|
|                      | B-1    | B-2  | B-3  | B-4  | METEOR      | ROUGE-L     | CIDER |
| یو و همکاران [۴۴]    | ۷۳,۱   | ۵۶,۵ | ۴۲,۴ | ۳۱,۶ | ۲۵,۰        | ۵۳,۳        | ۹۴,۳  |
| [۴۵] گو و همکاران    | ۶۷,۱   | ۴۷,۸ | ۳۲,۳ | ۲۱,۵ | ۲۰,۹        | ۴۷,۲        | ۶۹,۵  |
| شانکو و همکاران [۲۸] | -      | -    | -    | ۲۱,۹ | ۲۱,۱        | ۴۷,۲        | ۶۴    |
| لیان و همکاران [۴۶]  | ۶۱,۷   | ۴۲,۸ | ۲۸,۶ | ۱۹,۳ | ۲۰,۲        | ۴۵,۰        | ۶۱,۸  |
| فنگ و همکاران [۴۷]   | ۵۸,۹   | ۴۰,۳ | ۲۷,۰ | ۱۸,۶ | ۱۷,۹        | ۴۳,۱        | ۵۴,۹  |
| لو و همکاران [۱۹]    | ۷۴,۶   | ۵۸,۲ | ۴۴,۳ | ۳۳,۵ | ۲۶,۴        | ۵۵,۰        | ۱۰۶,۱ |
| چن و همکاران [۴۸]    | ۷۲,۵   | ۵۵,۶ | ۴۱,۴ | ۳۰,۶ | ۲۴,۶        | ۵۲,۸        | ۹۱,۱  |
| روش ما               | ۷۰,۶   | ۴۶,۱ | ۲۹,۱ | ۲۱,۲ | <u>۲۸,۴</u> | <u>۵۴,۶</u> | ۴۹,۵  |

Three dogs are playing in the snow



شکل (۷): توصیف تولید شده توسط مدل پیشنهادی



Ground-Truth: Dog looking out the window of a car in review mirror

MAGAN[31]: A dog looking out a car window in mirror

**Our-Method: a reflection of person holding a dog in the back of the car**

شکل (۸): مقایسه توصیف تولید شده با یکی از روش‌های پیشین

کردیم. لایه توجه چندگانه باعث می‌شود که مدل به قسمت‌های مهم تصویر بیشتر توجه کند و همچنین نمایش بهتری از تصویر نشان دهد. همچنین لایه توجه بر توجه باعث شد که مدل ما بتواند روابط بین اشیا را به خوبی درک کند. همچنان حوزه توصیف تصویر به صورت خودکار جای پیشرفت بسیاری دارد. همچنین چالش معیارهای ارزیابی همانطور که در مقدمه ذکر شد جای کار بسیار دارد و هنوز محققان در این حوزه با این مشکل رو به رو هستند.

## ۵- نتیجه‌گیری

با توجه به مزایای تولید توصیفات متنی تصویر توسط هوش مصنوعی، روش‌های بسیاری تا به امروز ارائه شده است. ما در این روش تلاش کردیم توصیفات تولید شده را بهبود بخشیم. همچنین مدل ما از منظر معیارهای METEOR و ROUGE دارای بهبود در عملکرد بود.

در این مقاله ما روشی مبتنی بر رمزگذار - رمزگشا پیشنهاد دادیم. و همچنین از لایه توجه چندگانه و توجه بر توجه استفاده

## مراجع

- Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3128-313 .<sup>۱۵</sup>
- [۱۵] A. Karpathy, A. Joulin, and L. F. Fei-Fei, "Deep fragment embeddings for bidirectional image sentence mapping," *Advances in neural information processing systems*, vol. 27, 2014.
- [۱۶] R. Kiros, R. Salakhutdinov, and R. S. Zemel, "Unifying visual-semantic embeddings with multimodal neural language models," *arXiv preprint arXiv:1411.2539*, 2014.
- [۱۷] Q. Wu, C. Shen, P. Wang, A. Dick, and A. Van Den Hengel, "Image captioning and visual question answering based on attributes and external knowledge," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 6, pp. 1367-1381, 2017.
- [۱۸] J. Li, N. Xu, W. Nie, and S. Zhang, "Image Captioning with multi-level similarity-guided semantic matching," *Visual Informatics*, vol. 5, no. 4, pp. 41-۴۸, ۲۰۲۱.
- [۱۹] J. Lu, C. Xiong, D. Parikh, and R. Socher, "Knowing when to look: Adaptive attention via a visual sentinel for image captioning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 375-383 .
- [۲۰] M. Z. Hossain, F. Sohel, M. F. Shiratuddin, and H. Laga, "A comprehensive survey of deep learning for image captioning," *ACM Computing Surveys (CSUR)*, vol. 51, no. 6, pp. 1-36, 2019.
- [۲۱] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.
- [۲۲] K. Cho, B. Van Merriënboer, D. Bahdanau, and Y. Bengio, "On the properties of neural machine translation: Encoder-decoder approaches," *arXiv preprint arXiv:1409.1259*, 2014.
- [۲۳] C. Wang and X. Gu, "Image captioning with adaptive incremental global context attention," *Applied Intelligence*, pp. 1-23, 2021.
- [۲۴] X. Chen and C. Lawrence Zitnick, "Mind's eye: A recurrent visual representation for image caption generation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 2422-2431 .
- [۲۵] K. Xu *et al.*, "Show, attend and tell: Neural image caption generation with visual attention," in *International conference on machine learning*, 2015: PMLR, pp. 2048-2057 .
- [۲۶] M. Pedersoli, T. Lucas, C. Schmid, and J. Verbeek, "Areas of attention for image captioning," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 1242-1250 .
- [۲۷] Z. Deng, Z. Jiang, R. Lan, W. Huang, and X. Luo, "Image captioning using DenseNet network and adaptive attention," *Signal Processing: Image Communication*, vol. 85, p. 115836, 2020.
- [۲۸] S. Cao, G. An, Z. Zheng, and Q. Ruan, "Interactions guided generative adversarial network for unsupervised image captioning," *Neurocomputing*, vol. 417, pp. 419-431, 2020.
- [۲۹] D. Wang, Z. Hu, Y. Zhou, R. Hong, and M. Wang, "A Text-Guided Generation and Refinement Model for
- [۱] J. Pavlopoulos, V. Kougia, and I. Androutsopoulos, "A survey on biomedical image captioning," in *Proceedings of the second workshop on shortcomings in vision and language*, 2019, pp. 26-36 .
- [۲] Z. F. Sattari, H. Khotanlou, and E. Alighardash, "Improving Image Captioning with Local Attention Mechanism," in *2022 9th Iranian Joint Congress on Fuzzy and Intelligent Systems (CFIS)*, 2022: IEEE, pp. 1-5 .
- [۳] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and tell: A neural image caption generator," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3156-3164 .
- [۴] X. Jia, E. Gavves, B. Fernando, and T. Tuytelaars, "Guiding the long-short term memory model for image caption generation," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2407-2415 .
- [۵] M. Zamiri and H. S. Yazdi, "Image annotation based on multi-view robust spectral clustering," *Journal of Visual Communication and Image Representation*, vol. 74, p. 103003, 2021.
- [۶] M. Zamiri, T. Bahraini, and H. S. Yazdi, "MVDF-RSC: Multi-view data fusion via robust spectral clustering for geo-tagged image tagging," *Expert Systems with Applications*, vol. 173, p. 114657, 2021.
- [۷] S. Li, G. Kulkarni, T. Berg, A. Berg, and Y. Choi , "Composing simple image descriptions using web-scale n-grams," in *Proceedings of the Fifteenth Conference on Computational Natural Language Learning*, 2011, pp. 220-228 .
- [۸] G. Kulkarni *et al.*, "Babytalk: Understanding and generating simple image descriptions," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 12, pp. 2891-2903, 2013.
- [۹] A. Farhadi *et al.*, "Every picture tells a story: Generating sentences from images," in *European conference on computer vision*, 2010: Springer ,pp. 15-29 .
- [۱۰] C. Sun, C. Gan, and R. Nevatia, "Automatic concept discovery from parallel text and visual corpora," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 2596-2604 .
- [۱۱] V. Ordonez, G. Kulkarni, and T. Berg " ,Im2text: Describing images using 1 million captioned photographs," *Advances in neural information processing systems*, vol. 24, 2011.
- [۱۲] Y. Gong, L. Wang, M. Hodosh, J. Hockenmaier, and S. Lazebnik, "Improving image-sentence embeddings using large weakly annotated photo collections," in *European conference on computer vision*, 2014: Springer, pp. 529-545 .
- [۱۳] M. Hodosh, P. Young, and J. Hockenmaier, "Framing image description as a ranking task: Data, models and evaluation metrics," *Journal of Artificial Intelligence Research*, vol. 47, pp. 853-899, 2013.
- [۱۴] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," in

captioning," *IEEE Transactions on Multimedia*, vol. 21, no. 11, pp. 2942-2956, 2019.

- [۴۴] Q. You, H. Jin, Z. Wang, C. Fang, and J. Luo, "Image captioning with semantic attention", in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4651-4659 .
- [۴۵] J. Gu, S. Joty, J. Cai, H. Zhao, X. Yang, and G. Wang, "Unpaired image captioning via scene graph alignments," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 10323-10332 .
- [۴۶] I. Laina, C. Rupprecht, and N. Navab, "Towards unsupervised image captioning with shared multimodal embeddings," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 7414-7424 .
- [۴۷] Y. Feng, L. Ma, W. Liu, and J. Luo, "Unsupervised image captioning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4125-4134 .
- [۴۸] L. Chen *et al.*, "Sca-cnn: Spatial and channel-wise attention in convolutional networks for image captioning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 5659-5667 .



زهرا فامیل ستاری در سال ۱۳۹۸ مدرک کارشناسی خود را در رشته مهندسی کامپیوتر از دانشگاه بوعلی سینا دریافت نمود و در سال ۱۴۰۰ مدرک کارشناسی ارشد خود را در رشته مهندسی کامپیوتر گرایش هوش مصنوعی از دانشگاه بوعلی سینا اخذ نمود. حوزه‌های پژوهشی مورد علاقه ایشان یادگیری عمیق، بینایی ماشین و پردازش زبان طبیعی است.



حسن ختن‌لو استاد گروه مهندسی کامپیوتر دانشگاه بوعلی سینا، در سال ۱۳۸۷ دکترای مهندسی کامپیوتر را از دانشگاه پیر و ماری کوری اخذ و تا به حال به عنوان عضو هیات علمی گروه مهندسی کامپیوتر، مشغول به فعالیت می‌باشد و زمینه‌های تحقیقاتی مورد علاقه ایشان، پردازش تصویر و ویدیو، پردازش تصاویر پزشکی، شناسایی الگو و یادگیری ماشین است.



الهام علیقارداش مربی گروه مهندسی کامپیوتر دانشگاه سید جمال الدین اسدآبادی، در سال ۱۳۹۱ مدرک کارشناسی ارشد خود را در رشته مهندسی کامپیوتر-گرایش هوش مصنوعی از دانشگاه بوعلی سینا اخذ نمود. وی در حال حاضر دانشجوی دکتری مهندسی کامپیوتر در دانشگاه بوعلی سینا می‌باشد.

زمینه‌های تحقیقاتی مورد علاقه ایشان بازشناسی الگو، پردازش تصویر، بینایی ماشین، یادگیری ماشین یادگیری چندوجهی و یادگیری عمیق است.

Image Captioning," *IEEE Transactions on Multimedia*, 2022.

- [۳۰] L. Huang, W. Wang, J. Chen, and X.-Y. Wei, "Attention on attention for image captioning," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 4634-4643 .
- [۳۱] Y. Wei, L. Wang, H. Cao, M. Shao, and C. Wu, "Multi-attention generative adversarial network for image captioning," *Neurocomputing*, vol. 387, pp. 91-99, 2020.
- [۳۲] T.-Y. Lin *et al.*, "Microsoft coco: Common objects in context," in *European conference on computer vision*, 2014: Springer, pp. 740-755 .
- [۳۳] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 2002, pp. 311-318 .
- [۳۴] S. Banerjee and A. Lavie, "METEOR: An automatic metric for MTEvaluation with improved correlation with human judgments," in *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 2005, pp. 65-72 .
- [۳۵] R. Vedantam, C. Lawrence Zitnick, and D. Parikh, "Cider: Consensus-based image description evaluation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 4566-4575 .
- [۳۶] M. Kilickaya, A. Erdem, N. Izkizler-Cinbis, and E. Erdem, "Re-evaluating automatic metrics for image captioning," *arXiv preprint arXiv:1612.07600*, 2016.
- [۳۷] C.-Y. Lin, "Rouge: A package for automatic evaluation of summaries," in *Text summarization branches out*, 2004, pp. 74-81 .
- [۳۸] K. Fu, J. Jin, R. Cui, F. Sha, and C. Zhang, "Aligning where to see and what to tell: Image captioning with region-based attention and scene-specific contexts," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 12, pp. 2321-2334, 2016.
- [۳۹] L. Li, S. Tang, Y. Zhang, L. Deng, and Q. Tian, "GLA: Global-local attention for image description," *IEEE Transactions on Multimedia*, vol. 20, no. 3, pp. 726-737, 2017.
- [۴۰] X. Zhu, L. Li, J. Liu, Z. Li, H. Peng, and X. Niu, "Image captioning with triple-attention and stack parallel LSTM," *Neurocomputing*, vol. 319, pp. 55-65, 2018.
- [۴۱] Y. Cheng, F. Huang, L. Zhou, C. Jin, Y. Zhang, and T. Zhang, "A hierarchical multimodal attention-based neural network for image captioning," in *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2017, pp. 889-892 .
- [۴۲] S. Ding, S. Qu, Y. Xi, and S. Wan, "Stimulus-driven and concept-driven analysis for image caption generation," *Neurocomputing*, vol. 398, pp. 520-530, 2020.
- [۴۳] X. Xiao, L. Wang, K. Ding, S. Xiang, and C. Pan, "Deep hierarchical encoder-decoder network for image