



دانشکده مهندسی کامپیوتر

## پرسش و پاسخ تصویری

گزارش سمینار کارشناسی ارشد در رشته مهندسی کامپیوتر  
گرایش هوش مصنوعی

مریم سادات هاشمی

استاد راهنما

دکتر سید صالح اعتمادی

دی ۱۳۹۹

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

## چکیده

مسئله پرسش و پاسخ تصویری یک مسئله چالش برانگیز است که در سال‌های اخیر معرفی شده است و مورد توجه بسیاری از محققان دو حوزه پردازش زبان طبیعی و بینایی ماشین قرار گرفته است. هدف این مسئله پاسخ به پرسش مطرح شده در مورد تصویر ورودی است. یک سیستم پرسش و پاسخ تصویری سعی می‌کند با استفاده از عناصر بصری تصویر و استنتاج جمع‌آوری شده از سوال متنی، پاسخ صحیح را پیدا کند. در فصل اول این بررسی، به معرفی مسئله پرسش و پاسخ تصویری، کاربرد و اهمیت آن و چالش‌های این مسئله می‌پردازیم. پس از تعریف برخی مفاهیم مورد نیاز در فصل دوم، مجموعه‌دادگان، روش‌های حل مسئله پرسش و پاسخ و تصویری و معیارهای ارزیابی آن را در فصل سوم بررسی می‌کنیم. با توجه به موفقیت یادگیری عمیق و مدل‌های از قبل آموزش دیده، رویکردهای حل مسئله پرسش و پاسخ تصویری را به دو دسته کلی رویکرد یادگیری عمیق و رویکرد مدل‌های از قبل آموزش دیده تقسیم‌بندی می‌کنیم. در فصل آخر، پس از نتیجه‌گیری در مورد ابعاد مختلف مسئله پرسش و پاسخ تصویری، در مورد مسیرهای تحقیق در آینده بحث می‌کنیم.

واژگان کلیدی: پرسش و پاسخ تصویری، پردازش زبان طبیعی، بینایی ماشین، یادگیری عمیق، مدل‌های از قبل آموزش دیده

# فهرست مطالب

چ	فهرست تصاویر
خ	فهرست جداول
۱	فصل ۱: مقدمه
۲	۱-۱ کاربرد و اهمیت مسئله
۳	۱-۲ بررسی چالش‌های موجود در این مسئله
۴	فصل ۲: تعاریف و مفاهیم مبنایی
۴	۲-۱ پردازش زبان طبیعی
۵	۲-۲ بینایی ماشین
۵	۲-۳ یادگیری عمیق
۶	۲-۴ شبکه‌های عصبی پیچشی
۷	۲-۴-۱ AlexNet
۷	۲-۴-۲ VGGNet
۸	۲-۴-۳ GoogleNet
۹	۲-۴-۴ ResNet
۹	۲-۵ شبکه‌های عصبی بازگشتی
۹	۲-۵-۱ LSTM
۱۰	۲-۵-۲ GRU
۱۰	۲-۶ تعبیه کلمات

۱۱	۲-۶-۱ کدگذاری one-hot
۱۱	۲-۶-۲ Skip-gram و CBOW
۱۱	۲-۶-۳ GloVe
۱۲	۲-۶-۴ LSTM، CNN و GRU
۱۲	۲-۷ جمع‌بندی
۱۳	فصل ۳: مروری بر کارهای مرتبط
۱۳	۳-۱ بررسی مجموعه داده‌گان مطرح این حوزه
۱۴	۳-۱-۱ داده‌گان DAQUAR
۱۴	۳-۱-۲ داده‌گان VQA
۱۶	۳-۱-۳ Visual Madlibs
۱۷	۳-۱-۴ Visual7w
۱۹	۳-۱-۵ CLEVR
۱۹	۳-۱-۶ Tally-QA
۲۰	۳-۱-۷ KVQA
۲۱	۳-۲ تقویت داده‌گان در مسئله پرسش و پاسخ تصویری
۲۲	۳-۳ رویکرد یادگیری عمیق
۲۳	۳-۳-۱ فاز ۱: استخراج ویژگی از تصویر و سوال
۲۶	۳-۳-۲ فاز ۲: بازنمایی مشترک تصویر و سوال
۲۶	۳-۳-۲-۱ روش‌های پایه
۲۶	۳-۳-۲-۲ روش‌های مبتنی بر شبکه‌های عصبی
۲۹	۳-۳-۲-۳ روش‌های مبتنی بر مکانیزم توجه
۳۰	۳-۳-۳ فاز ۳: پیش‌بینی پاسخ
۳۱	۳-۴ رویکرد مدل‌های از قبل آموزش‌دیده بر روی زبان طبیعی و تصویر
۳۱	۳-۴-۱ معماری تک جریان
۳۲	۳-۴-۱-۱ شبکه VL-BERT

۳۴	۳-۴-۱-۲ شبکه UNITER
۳۵	۳-۴-۱-۳ شبکه VLP
۳۶	۳-۴-۱-۴ شبکه OSCAR
۳۷	۳-۴-۲ معماری دو جریان
۳۷	۳-۴-۲-۱ شبکه ViLBERT
۳۹	۳-۴-۲-۲ شبکه LXMERT
۴۱	۳-۵ معیارهای ارزیابی مسئله پرسش و پاسخ تصویری
۴۲	۳-۵-۱ معیار دقت
۴۲	۳-۵-۲ معیار شباهت Wu-Palmer
۴۲	۳-۵-۳ معیار اجماع
۴۳	۳-۵-۴ معیار MPT
۴۳	۳-۵-۵ معیار BLEU
۴۳	۳-۵-۶ معیار METEOR
۴۴	۳-۶ جمع بندی
۴۵	فصل ۴: نتیجه گیری و کارهای آینده
۴۵	۴-۱ نتیجه گیری
۴۶	۴-۲ مسائل باز و کارهای قابل انجام
۴۸	مراجع
۵۶	واژه نامه فارسی به انگلیسی
۵۸	واژه نامه انگلیسی به فارسی

## فهرست تصاویر

- ۱-۱ مثالی از سیستم پرسش و پاسخ متنی و تصویری . . . . . ۲
- ۱-۲ نمونه‌ای از شبکه عصبی عمیق . . . . . ۶
- ۲-۲ معماری شبکه AlexNet . . . . . ۷
- ۳-۲ معماری شبکه VGGNet . . . . . ۷
- ۴-۲ معماری شبکه GoogleNet . . . . . ۸
- ۵-۲ معماری شبکه ResNet . . . . . ۸
- ۶-۲ مقایسه معماری شبکه‌های عصبی بازگشتی، LSTM و GRU . . . . . ۱۰
- ۷-۲ Skip-gram و CBOW . . . . . ۱۲
- ۱-۳ چند نمونه از دادگان DAQUAR . . . . . ۱۵
- ۲-۳ چند نمونه از دادگان VQA v1 - real . . . . . ۱۶
- ۳-۳ چند نمونه از دادگان VQA v1 - abstarct . . . . . ۱۶
- ۴-۳ چند نمونه از دادگان VQA v2 . . . . . ۱۷
- ۵-۳ یک نمونه از دادگان Visual Madlibs . . . . . ۱۸
- ۶-۳ چند نمونه از دادگان Visual7W . . . . . ۱۸
- ۷-۳ چند نمونه از دادگان CLEVR . . . . . ۱۹
- ۸-۳ چند نمونه از دادگان Tally-QA . . . . . ۲۰
- ۹-۳ چند نمونه از دادگان KVQA . . . . . ۲۱
- ۱۰-۳ حالت اول معماری رمزگذار-رمزگشا در پرسش و پاسخ تصویری . . . . . ۲۸
- ۱۱-۳ حالت دوم معماری رمزگذار-رمزگشا در پرسش و پاسخ تصویری . . . . . ۲۸

۳-۱۲ معماری شبکه از قبل آموزش دیده VL-BERT . . . . .	۳۲
۳-۱۳ نحوه ورودی و خروجی شبکه VL-BERT برای آموزش در مسئله پرسش و پاسخ تصویری . . . . .	۳۳
۳-۱۴ معماری شبکه از قبل آموزش دیده UNITER . . . . .	۳۴
۳-۱۵ معماری شبکه از قبل آموزش دیده VLP . . . . .	۳۵
۳-۱۶ معماری شبکه از قبل آموزش دیده OSCAR . . . . .	۳۷
۳-۱۷ معماری شبکه از قبل آموزش دیده ViLBERT . . . . .	۳۷
۳-۱۸ ساختار لایه co-attentional transformer . . . . .	۳۸
۳-۱۹ معماری شبکه از قبل آموزش دیده LXMERT . . . . .	۳۹



## فهرست جداول

- ۲-۱ مقایسه مهم‌ترین شبکه‌های عصبی پیچشی آموزش دیده بر روی دادگان ImageNet . . . . ۷
- ۳-۱ بررسی مجموعه‌دادگان در حوزه پرسش و پاسخ تصویری. . . . . ۱۴
- ۳-۲ الگوهای استفاده شده برای تولید سوال در دادگان DAQUAR. . . . . ۱۵
- ۳-۳ شبکه‌های عصبی پیچشی استفاده شده در مدل‌های پرسش و پاسخ تصویری. . . . . ۲۴
- ۳-۴ تعبیه کلمات استفاده شده در مدل‌های پرسش و پاسخ تصویری. . . . . ۲۵
- ۳-۵ بررسی رویکرد پیش‌بینی پاسخ در چند نمونه از مدل‌های پرسش و پاسخ تصویری. . . . . ۳۰
- ۳-۶ مقایسه بین شبکه‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر . . . . . ۴۰
- ۳-۷ دقت شبکه‌های از قبل آموزش دیده بر روی مجموعه‌داده VQA v2.0 (test-std) . . . . . ۴۱

# فصل ۱

## مقدمه

در سال‌های اخیر پیشرفت‌های زیادی در مسائل هوش مصنوعی و یادگیری عمیق که در تقاطع دو حوزه پردازش زبان طبیعی و بینایی ماشین قرار می‌گیرند؛ رخ داده است. یکی از مسائلی که اخیراً مورد توجه قرار گرفته است؛ پرسش و پاسخ تصویری است. با توجه به یک تصویر و یک سؤال به زبان طبیعی، سیستم سعی می‌کند با استفاده از عناصر بصری تصویر و استنتاج جمع‌آوری شده از سوال متنی، پاسخ صحیح را پیدا کند [۴۴]. پرسش و پاسخ تصویری نسخه گسترش‌یافته مسئله پرسش و پاسخ متنی است که اطلاعات بصری به مسئله اضافه شده است. شکل ۱-۱ گویای تفاوت این دو مسئله است.

در سیستم پرسش و پاسخ متنی، یک متن و یک سوال متنی به عنوان ورودی به سیستم داده می‌شود و انتظار می‌رود که سیستم با توجه به درک و تفسیری که از متن و سوال بدست می‌آورد؛ یک جواب متنی را خروجی دهد. اما در سیستم پرسش و پاسخ تصویری، یک تصویر و یک سوال متنی به ورودی سیستم داده می‌شود و انتظار می‌رود که سیستم بتواند با استفاده از عناصر بصری تصویر و تفسیری که از سوال بدست می‌آورد؛ یک پاسخ متنی را در خروجی نشان دهد.

مسئله پرسش و پاسخ تصویری پیچیدگی بیشتری نسبت به مسئله پرسش و پاسخ متنی دارد زیرا تصاویر بعد بالاتر و نویز بیشتری نسبت به متن دارند. علاوه بر این، تصاویر فاقد ساختار و قواعد دستوری زبان هستند. در نهایت هم، تصاویر غنای بیشتری از دنیای واقعی را ضبط می‌کنند، در حالی که زبان طبیعی در حال حاضر نشانگر سطح بالاتری از انتزاع دنیای واقعی است [۷۲].



شکل ۱-۱: مثالی از سیستم پرسش و پاسخ متنی و تصویری

## ۱-۱ کاربرد و اهمیت مسئله

در طی سال‌های متمادی، محققان به دنبال ساخت ماشین‌هایی بودند که به اندازه‌ی کافی هوشمند باشند که از آن به طور موثر همانند انسان‌ها برای تعامل استفاده کنند. مسئله‌ی پرسش و پاسخ تصویری یکی از پله‌های رسیدن به این رویای هوش مصنوعی است و از این جهت حائز اهمیت است.

کاربردهای بسیاری برای پرسش و پاسخ تصویری وجود دارد. یکی از مهم‌ترین موارد دستیار هوشمند برای افراد کم‌بینا و نابینا است [۲۰]. علاوه بر این، در سال‌های اخیر دستیاران صوتی<sup>۱</sup> و عامل‌های گفتگو<sup>۲</sup> مانند Cortana، Siri و Alexa در بازار عرضه شدند که می‌توانند با انسان‌ها با استفاده از زبان طبیعی ارتباط برقرار کنند. در حال حاضر این دستیاران با استفاده از صوت و متن این ارتباط را برقرار می‌کنند در نتیجه گفتگوی بین این دستیاران با انسان‌ها مشابه دنیای واقعی نمی‌باشد. این ارتباط را می‌توان با استفاده از داده‌های تصویری و ویدئویی به واقعیت نزدیک‌تر کرد. اینجاست که مسئله‌ی پرسش و پاسخ تصویری برای نزدیک کردن تعامل بین انسان و عامل‌های گفتگو به دنیای واقعی می‌تواند موثر باشد. همین موضوع را می‌توانیم به صورت گسترده‌تری در ربات‌ها مشاهده کنیم. برای این‌که ربات بتواند بهتر با انسان‌ها ارتباط برقرار کند و به سوالات و درخواست‌ها پاسخ دهد؛ نیاز دارد که درک و فهم درستی از اطراف داشته باشد که این مستلزم داشتن تصویری دقیق از پیرامون است. بنابراین این ربات می‌تواند برای پاسخ به پرسش‌ها از

<sup>۱</sup> voice assistants  
<sup>۲</sup> conversational agents

دانشی که از طریق تصویر پیرامون خود بدست می‌آورد، جواب درستی را بدهد. کاربرد دیگر این مسئله در پزشکی است. در بسیاری از موارد تحلیل تصاویر پزشکی مانند تصاویر CT اسکن و x-ray برای یک پزشک متخصص هم دشوار است. اما یک سیستم پرسش و پاسخ تصویری می‌تواند با تحلیل و تشخیص موارد غیرطبیعی موجود در تصویر، به عنوان نظر دوم به پزشک متخصص کمک کند. از طرفی ممکن است در بعضی اوقات بیمار دسترسی به پزشک را نداشته باشد تا شرح تصاویر را متوجه شود. وجود سیستم پرسش و پاسخ تصویری می‌تواند آگاهی بیمار را نسبت به بیماری افزایش دهد و از نگرانی او بکاهد [۶۴].

## ۱-۲. بررسی چالش‌های موجود در این مسئله

در مقایسه با مسائل دیگری که مشترک بین پردازش زبان طبیعی و بینایی ماشین است مانند توصیف تصویر<sup>۳</sup> و بازیابی متن به تصویر<sup>۴</sup>، مسئله پرسش و پاسخ تصویری چالش‌برانگیزتر است زیرا (۱) سوالات از پیش تعیین نشده است. به این معنی که در مسئله‌ای مانند تشخیص اشیا، سوال این است که چه اشیایی در تصویر وجود دارد و این سوال از پیش تعیین شده است و در طول حل مسئله تغییر نمی‌کند و تنها تصویر تغییر می‌کند که منجر به پاسخ‌های متفاوت می‌شود. اما در پرسش و پاسخ تصویری، برای هر تصویر سوالات متفاوت و مرتبط با همان تصویر پرسیده می‌شود که در زمان اجرا تعیین می‌شود. (۲) اطلاعات موجود در تصویر ابعاد بالایی دارد که پردازش آن‌ها به زمان و حافظه زیادی نیاز دارد. (۳) مسئله پرسش و پاسخ تصویری نیاز به حل مسائل پایه‌ای و فرعی دارد مانند تشخیص اشیا<sup>۵</sup> (آیا در تصویر سگ وجود دارد؟)، تشخیص فعالیت<sup>۶</sup> (آیا کودک گریه می‌کند؟)، طبقه‌بندی صفات<sup>۷</sup> (چتر چه رنگی است؟)، شمارش (چند نفر در تصویر وجود دارد؟)، طبقه‌بندی صحنه<sup>۸</sup> (هوا بارانی است؟) و روابط مکانی بین اشیا (چه چیزی بین گربه و مبل است؟).

<sup>۳</sup> image captioning  
<sup>۴</sup> text-to-image retrieval  
<sup>۵</sup> object detection  
<sup>۶</sup> activity recognition  
<sup>۷</sup> attribute classification  
<sup>۸</sup> scene classification

## فصل ۲

### تعاریف و مفاهیم مبنایی

همان‌طور که قبلاً اشاره شد، مسئله پرسش و پاسخ تصویری در تقاطع دو حوزه پردازش زبان طبیعی و بینایی ماشین قرار می‌گیرد. از این رو قبل از بررسی کارهای مرتبط با مسئله پرسش و پاسخ تصویری، نیاز است تا با مفاهیم مربوط به این دو حوزه آشنا شویم. در ادامه این فصل به شرح مفاهیم و تعاریف پایه می‌پردازیم.

#### ۲-۱ پردازش زبان طبیعی

پردازش زبان طبیعی<sup>۱</sup> یکی از زیرشاخه‌های علوم کامپیوتر و هوش مصنوعی است که به تعامل بین کامپیوتر و زبان‌های (طبیعی) انسانی می‌پردازد. هدف اصلی در پردازش زبان طبیعی، تحلیل زبان‌های طبیعی به منظور آسان‌تر ساختن فهم آن‌ها برای کامپیوتر می‌باشد. مسلماً در صورتی که کامپیوتر بتواند توسط زبان‌های طبیعی با انسان ارتباط برقرار کند، بسیاری از مشکلات تعامل انسان با کامپیوتر حل شده و زندگی برای انسان‌ها راحت‌تر خواهد شد. با پیشرفت تکنولوژی و بوجود آمدن نیازهای متفاوت برای انسان‌ها، کاربردهای جدیدی برای این حوزه تعریف می‌شود. ترجمه ماشینی<sup>۲</sup>، خلاصه‌سازی متون<sup>۳</sup>، تحلیل احساسات<sup>۴</sup>، طبقه‌بندی متون<sup>۵</sup>،

---

<sup>۱</sup> natural language processing

<sup>۲</sup> machine translation

<sup>۳</sup> text summarization

<sup>۴</sup> sentiment analysis

<sup>۵</sup> text classification

سیستم‌های توصیه‌گر<sup>۶</sup>، غلطیابی متون<sup>۷</sup> از جمله مهم‌ترین کاربردهای پردازش زبان طبیعی است.

## ۲-۲ بینایی ماشین

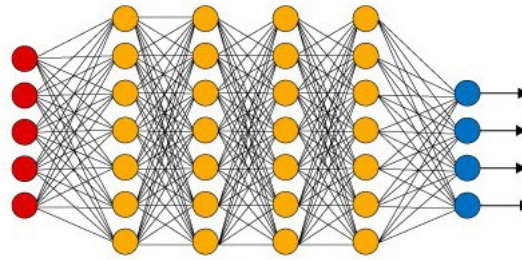
بینایی ماشین<sup>۸</sup> جزو حوزه‌های در حال توسعه در علوم کامپیوتر و هوش مصنوعی محسوب می‌شود که سعی دارد از طریق پردازش تصاویر دوبعدی، جهان سه‌بعدی پیرامون را بازسازی و تفسیر کند. به بیان ساده، بینایی ماشین یعنی کامپیوترها بتوانند جهان را به کمک دوربین‌ها ببینند، بفهمند و حتی از بینایی انسان پیشی بگیرند. بینایی کامپیوتر دارای کاربردهای بسیار متنوعی مانند طبقه‌بندی اشیا<sup>۹</sup>، تشخیص اشیا<sup>۱۰</sup>، تقسیم‌بندی اشیا<sup>۱۱</sup>، تشخیص چهره<sup>۱۲</sup> است.

## ۳-۲ یادگیری عمیق

یادگیری عمیق<sup>۱۳</sup> زیر شاخه‌ای از یادگیری ماشین<sup>۱۴</sup> است که تلاش می‌کند مفاهیم انتزاعی سطح بالا را با استفاده از نمونه‌های (دادگان) زیاد مدل نماید. بیشتر روش‌های یادگیری عمیق از معماری شبکه‌های عصبی مصنوعی<sup>۱۵</sup> استفاده می‌کنند. به همین دلیل است که اغلب از مدل‌های یادگیری عمیق به عنوان شبکه‌های عصبی عمیق یاد می‌شود. اصطلاح «عمیق» معمولاً به تعداد لایه‌های پنهان در شبکه عصبی اشاره دارد. شبکه‌های عصبی سستی فقط شامل ۲ یا ۳ لایه پنهان هستند، در حالی که شبکه‌های عمیق می‌توانند تا ۱۵۰ لایه داشته باشند. مدل‌های یادگیری عمیق معمولاً با استفاده از مجموعه‌های بزرگی از داده‌های دارای برچسب آموزش می‌بینند. نکته‌ی حائز اهمیت در شبکه‌های عصبی عمیق یادگیری مستقیم ویژگی‌ها بدون نیاز به استخراج دستی می‌باشد.

در حالی که یادگیری عمیق برای اولین بار در دهه ۱۹۸۰ مطرح شد اما به دلیل تولید داده‌های زیاد، افزایش

recommender systems<sup>۶</sup>  
 spell correction<sup>۷</sup>  
 computer vision<sup>۸</sup>  
 object classification<sup>۹</sup>  
 object detection<sup>۱۰</sup>  
 object segmentation<sup>۱۱</sup>  
 face recognition<sup>۱۲</sup>  
 deep learning<sup>۱۳</sup>  
 machine learning<sup>۱۴</sup>  
 artificial neural networks<sup>۱۵</sup>



شکل ۲-۱: نمونه‌ای از شبکه عصبی عمیق. نوروهای قرمز لایه ورودی، نوروهای نارنجی لایه مخفی و نوروهای آبی لایه خروجی را نشان می‌دهند.

قدرت محاسباتی و پیشرفت الگوریتم‌های این حوزه شاهد پیشرفت چشمگیر یادگیری عمیق در سال‌های اخیر هستیم. در حال حاضر شبکه‌های عصبی عمیق در حوزه‌های زیادی از جمله پردازش زبان طبیعی، بینایی ماشین، پردازش گفتار کاربرد دارد. شبکه‌های عصبی پیچشی و شبکه‌های عصبی بازگشتی از مهم‌ترین و پرکاربردترین شبکه‌های یادگیری عمیق هستند.

## ۲-۴ شبکه‌های عصبی پیچشی

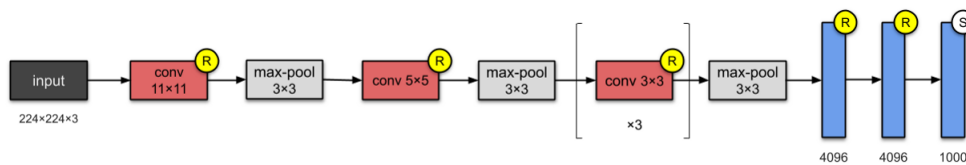
شبکه‌های عصبی پیچشی<sup>۱۶</sup> دسته‌ای از شبکه‌های عصبی عمیق هستند که معمولاً برای تجزیه و تحلیل تصاویر استفاده می‌شوند. برخلاف شبکه‌های کاملاً متصل<sup>۱۷</sup> که هر نورون در یک لایه به همه نورون‌های لایه بعدی متصل است، در شبکه‌های عصبی پیچشی هر نورون تنها به بخشی از نورون‌های لایه بعدی متصل است. این خاصیت به دلیل انجام عملیات کانولوشن در شبکه‌های عصبی پیچشی است و باعث می‌شود که الگوهای محلی را از داده فرا بگیرند. در حالی که شبکه‌های کاملاً متصل الگوهای سراسری را یاد می‌گیرند. معمولاً از شبکه‌های عصبی پیچشی برای استخراج ویژگی از تصویر استفاده می‌شود. در ادامه چند نمونه از برجسته‌ترین شبکه‌های عصبی پیچشی را معرفی می‌کنیم.

<sup>۱۶</sup>convolutional neural networks

<sup>۱۷</sup>fully connected networks

جدول ۲-۱: مقایسه مهم‌ترین شبکه‌های عصبی پیچشی آموزش دیده بر روی دادگان ImageNet [۱۳]

مدل CNN	سال	تعداد لایه‌ها	ابعاد ورودی	ابعاد خروجی	پارامترها
AlexNet [۲۲]	۲۰۱۲	۸	$227 \times 227$	۴۰۹۶	۶۰ میلیون
VGGNet [۶۰]	۲۰۱۴	۱۹	$224 \times 224$	۴۰۹۶	۱۳۸ میلیون
GoogleNet [۶۳]	۲۰۱۴	۲۲	$229 \times 229$	۱۰۲۴	۵ میلیون
ResNet [۲۱]	۲۰۱۵	۱۵۲	$224 \times 224$	۲۰۱۴۸	۲۶ میلیون



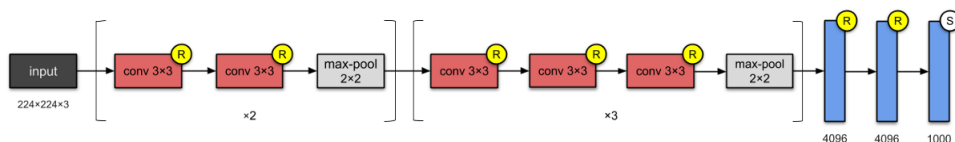
شکل ۲-۲: معماری شبکه AlexNet

## ۲-۴-۱ AlexNet

در سال ۲۰۱۲، شبکه AlexNet [۲۲] ارائه شد. در این شبکه از ۵ لایه پیچشی و ۳ لایه کاملاً متصل استفاده شده است. معماری این شبکه در شکل ۲-۲ نمایش داده شده است.

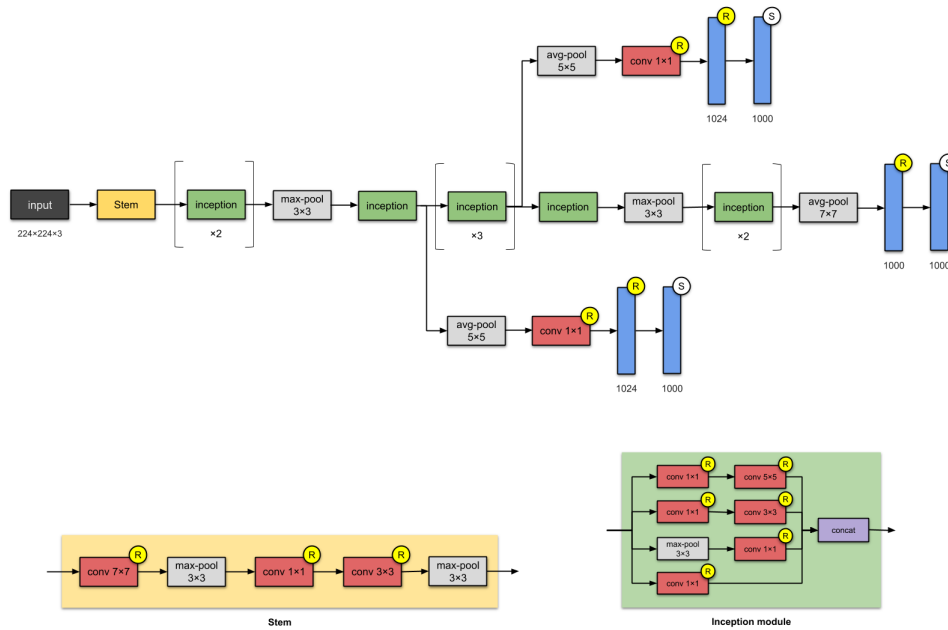
## ۲-۴-۲ VGGNet

شبکه VGGNet [۶۰] در سال ۲۰۱۴ معرفی شد. شبکه VGGNet از ۱۶ لایه پیچشی تشکیل شده است و معماری بسیار یکنواختی دارد. شبکه VGGNet یکی از محبوب‌ترین شبکه‌ها برای استخراج ویژگی است. تعداد پارامترهای این شبکه برابر با ۱۳۸ میلیون است. معماری این شبکه در شکل ۲-۳ نشان داده شده است.

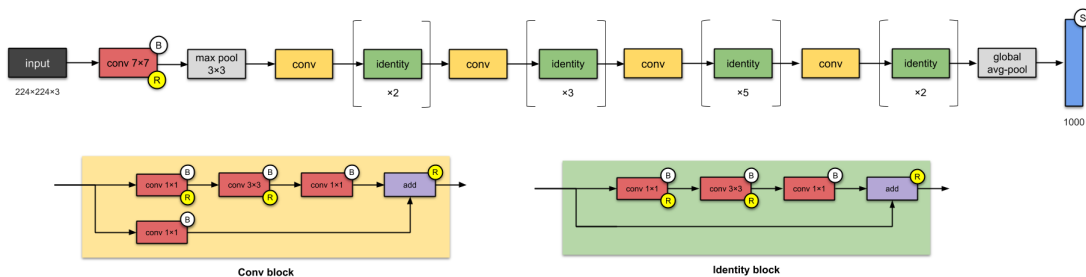


شکل ۲-۳: معماری شبکه VGGNet





شکل ۲-۴: معماری شبکه GoogleNet



شکل ۲-۵: معماری شبکه ResNet

## ۲-۴-۳ GoogleNet

شبکه GoogleNet [۶۳] همزمان با شبکه VGGNet در سال ۲۰۱۴ معرفی شد و توانست بر شبکه VGGNet غلبه کند و دقت بالاتری را بر روی دادگان ImageNet بدست آورد. مهم‌ترین علت موفقیت آن استفاده از ماژول inception بود که منجر به کاهش شدید تعداد پارامترها در این شبکه شد. شبکه GoogleNet از ۲۲ لایه با ۴ میلیون پارامتر تشکیل شده است. معماری این شبکه در شکل ۲-۴ نمایش داده شده است.

## ResNet ۴-۴-۲

شبکه ResNet [۲۱] با معرفی مفهوم جدید skip connection این امکان را برای شبکه‌های عصبی پیچشی ایجاد کرد که شبکه‌ها عمیق‌تر شوند و در عین حال آموزش در زمان کمتری انجام شود. با وجود اتصالات skip connection ورودی هر لایه بدون واسطه به لایه بعدی منتقل می‌شود. بنابراین مشکل از محوشدگی گرادیان<sup>۱۸</sup> در شبکه‌های عمیق رفع می‌شود. ۱۵۲ لایه در شبکه ResNet به کار رفته است. معماری این شبکه در شکل ۲-۵ نشان داده شده است.

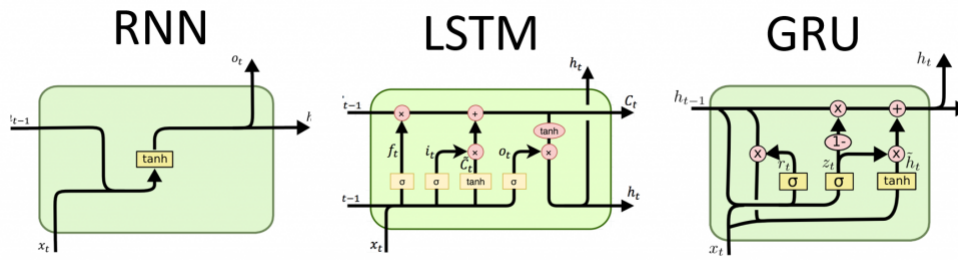
## ۲-۵ شبکه‌های عصبی بازگشتی

شبکه‌های عصبی بازگشتی<sup>۱۹</sup> دسته‌ی دیگری از شبکه‌های عصبی عمیق هستند که معمولاً برای پردازش داده‌های دنباله‌دار مانند جملات، صوت و ویدئو استفاده می‌شوند. این شبکه‌ها دارای یک نوع حافظه هستند که اطلاعاتی که تا کنون دیده‌اند را ضبط می‌کنند. در تئوری، شبکه‌های عصبی بازگشتی می‌توانند اطلاعات موجود در یک دنباله طولانی را ضبط و از آن‌ها استفاده کنند اما در عمل این‌طور نیست و بسیار محدود هستند و فقط اطلاعات چند گام قبل را نگه می‌دارند. شبکه‌های عصبی بازگشتی پارامترهای مشابهی را بین همه گام‌های زمانی به اشتراک می‌گذارند. این بدین معنی است که در هر گام زمانی عملیات مشابهی را انجام می‌دهند و فقط ورودی‌ها متفاوت هستند. با این تکنیک تعداد کلی پارامترهایی که شبکه باید یاد بگیرد به شدت کاهش پیدا می‌کند. در ادامه این بخش به معرفی دو شبکه عصبی بازگشتی معروف می‌پردازیم.

## LSTM ۱-۵-۲

مدل LSTM<sup>۲۰</sup> [۲۳] در سال ۱۹۹۵ برای توسعه شبکه‌های عصبی بازگشتی ارائه شد. شبکه LSTM برای حل مشکل محوشدگی گرادیان در شبکه‌های عصبی بازگشتی بوجود آمد. بزرگ‌ترین ویژگی LSTM امکان یادگیری وابستگی بلند مدت است که توسط شبکه‌های عصبی بازگشتی امکان‌پذیر نبود. برای پیش‌بینی گام زمانی بعدی نیاز است که مقادیر وزن‌ها در شبکه بروزرسانی شوند که این کار مستلزم حفظ اطلاعات گام‌های زمانی ابتدایی است. یک شبکه عصبی بازگشتی فقط می‌تواند تعداد محدودی از وابستگی‌های کوتاه مدت را

<sup>۱۸</sup> vanishing gradient<sup>۱۹</sup> recurrent neural networks<sup>۲۰</sup> Long Short-Term Memory



شکل ۶-۲: مقایسه معماری شبکه‌های عصبی بازگشتی، LSTM و GRU [منبع]

یاد بگیرد و نمی‌تواند سری‌های زمانی بلندمدت را یاد بگیرد. اما LSTM می‌تواند این وابستگی‌های بلند مدت را به درستی یاد بگیرد.

## ۲-۵-۲ GRU

یکی دیگر از شبکه‌های عصبی بازگشتی، GRU<sup>۲۱</sup> [۱۱] است که در سال ۲۰۱۴ معرفی شد. این شبکه نیز مانند LSTM مشکل محوشدگی گرادیان در شبکه‌های عصبی بازگشتی را حل می‌کند. در واقع GRU نوع خاصی از LSTM است که با کم کردن تعداد دروازه‌ها، سرعت محاسبات را افزایش داده است.

## ۶-۲ تعبیه کلمات

بیشتر الگوریتم‌های یادگیری ماشین و یادگیری عمیق قادر به پردازش متن به شکل خام و ساده نیستند و برای بازنمایی متن‌ها نیاز به تعبیه کلمات<sup>۲۲</sup> دارند. تعبیه کلمات نگاشت کلمات یا عبارات از واژگان به بردارهای عددی است تا کامپیوترها بتوانند به راحتی آن‌ها را پردازش کنند. تعبیه کلمات عمدتاً برای مدل‌سازی زبان و یادگیری ویژگی در پردازش زبان طبیعی استفاده می‌شود. ایده اصلی در پشت تمام روش‌های تعبیه کلمات، گرفتن هرچه بیشتر اطلاعات معنایی و ریخت‌شناسی است. روش‌های تعبیه کلمات بسیاری در مسئله پرسش و پاسخ تصویری استفاده شده است. در ادامه به برجسته‌ترین و پرکاربردترین روش‌های تعبیه کلمات موجود و استفاده‌شده در مسئله پرسش و پاسخ تصویری می‌پردازیم.

<sup>۲۱</sup>Gated Recurrent Unit  
<sup>۲۲</sup>word embedding

## ۲-۶-۱ کدگذاری one-hot

روش کدگذاری one-hot ساده‌ترین روش تعبیه کلمات است. در این روش یک لغت‌نامه از همه واژه‌های منحصر به فرد موجود در دادگان ساخته می‌شود و اندیس یکتایی به هر واژه اختصاص می‌یابد. بنابراین برای هر واژه یک بردار به طول تعداد واژه‌ها ساخته می‌شود که تمامی مقادیر آن صفر است به جز اندیس مربوط به همان واژه که مقدار آن یک است. پیاده‌سازی این روش آسان است اما طول بردارها بزرگ است زیرا برابر با تعداد کل واژه‌های منحصر به فرد دادگان است و هزینه زیادی برای ذخیره‌سازی دارد. بزرگ‌ترین عیب این روش این است که نمی‌توان از آن معنا و مفهوم استخراج کرد زیرا فاصله‌ی تمامی کلمات با هم یکسان است. در صورتی که ما انتظار داریم؛ کلماتی که مشابه هم هستند بردارهای نزدیک به هم یا مشابه هم داشته باشند و کلماتی که معنای متفاوتی با یکدیگر دارند تا حد امکان بردارهایشان از هم دور باشند.

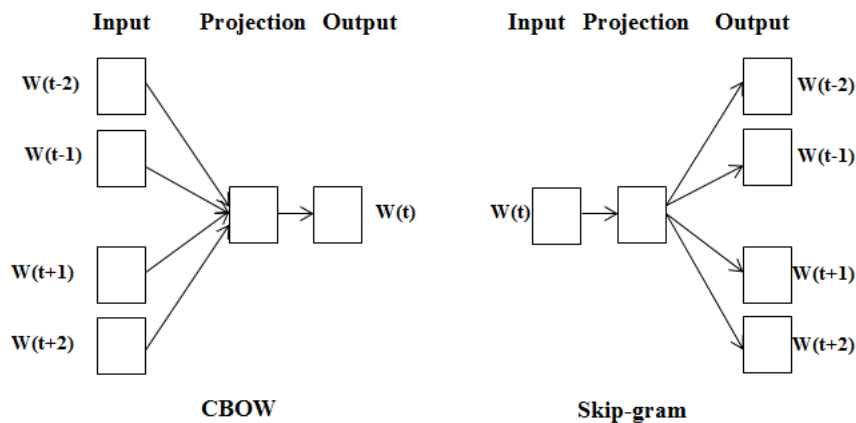
## ۲-۶-۲ CBOW و Skip-gram

برای رفع مشکلات کدگذاری one-hot، دو روش CBOW<sup>۲۳</sup> [۴۵] و Skip-gram [۴۵] پیشنهاد شد که از شبکه‌های عصبی به عنوان جز اصلی خود استفاده می‌کنند. این دو مدل بر عکس هم کار می‌کنند. در هر دو مدل، از یک شبکه عصبی سه لایه که شامل لایه ورودی، لایه پنهان و لایه خروجی است، استفاده شده است. در مدل CBOW کلمات اطراف و نزدیک به یک کلمه (n-1 کلمه) به لایه ورودی داده می‌شود و مدل سعی می‌کند این کلمه (nامین کلمه) را حدس بزند. بعد از آموزش این شبکه، وزن بین لایه‌ی پنهان و لایه خروجی، کلمات دادگان را بازنمایی می‌کند که هر ستون آن بردار مربوط به یک کلمه را نشان می‌دهد. در مدل skip-gram برعکس CBOW یک کلمه به شبکه ورودی داده می‌شود و شبکه باید کلمات اطراف و نزدیک به آن را حدس بزند. معماری CBOW و Skip-gram در شکل ۲-۷ آورده شده است.

## ۲-۶-۳ GloVe

یکی دیگر از تعبیه کلمات مشهور، مدل بردار سراسری یا به اختصار GloVe<sup>۲۴</sup> است که توسط پنینگتون و همکاران [۴۹] در سال ۲۰۱۴ در تیم پردازش زبان‌های طبیعی دانشگاه استنفورد معرفی و توسعه داده شد. در GloVe فاصله میان بردارها نشان‌دهنده شباهت معنایی میان آن بردارها است.

<sup>۲۳</sup> Continouse Bag Of Words<sup>۲۴</sup> Global Vector



شکل ۷-۲: معماری شبکه CBOW و Skip-gram [۴۵]

## ۴-۶-۲ GRU و LSTM ، CNN

با پیشرفت یادگیری عمیق در دهه اخیر، محققان برای استخراج ویژگی و بازنمایی متن از LSTM ، CNN ، [۲۳] و GRU [۱۱] استفاده می‌کنند. برای استخراج ویژگی از متن با استفاده از CNN بردارهای کلمات در کنار هم قرار داده می‌شود سپس به لایه‌های کانولوشنی یک بعدی داده می‌شود و فیلترهای متفاوتی بر روی آن‌ها اعمال می‌شود و پس از عبور از لایه max-pooling ویژگی‌ها بدست می‌آید. همچنین برای استخراج ویژگی از متن با استفاده از LSTM و GRU کافی است، بردار کلمات یک جمله به عنوان ورودی به این لایه‌ها داده شود. سپس خروجی آخرین گام زمانی به عنوان ویژگی کل جمله خواهد بود.

## ۷-۲ جمع‌بندی

در این فصل، دو حوزه پردازش زبان طبیعی و بینایی ماشین معرفی گردید. پس از آن با یادگیری عمیق به عنوان یکی از زیرشاخه‌های یادگیری ماشین آشنا شدیم. در ادامه به بررسی شبکه‌های عصبی پیچشی و شبکه‌های عصبی بازگشتی به عنوان مهم‌ترین شبکه‌های عصبی عمیق پرداختیم. در انتها آموختیم که چگونه از این دو برای استخراج ویژگی از تصویر و متن استفاده می‌کنیم. آشنایی با این مفاهیم به درک راه‌حل‌های پیشنهاد شده در مسئله پرسش و پاسخ تصویری به ما کمک خواهد کرد.

## فصل ۳

### مروری بر کارهای مرتبط

بسیاری از محققان راه‌حل‌ها یا الگوریتم‌هایی را برای حل مسئله پرسش و پاسخ تصویری پیشنهاد کرده‌اند. در این بخش ما آن‌ها را به دو رویکرد کلی تقسیم می‌کنیم: رویکرد یادگیری عمیق و رویکرد شبکه‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر. از آنجایی که برای آموزش مدل‌های یادگیری عمیق نیاز به مجموعه‌داده‌گان است، در ابتدا، به بررسی و معرفی مجموعه‌داده‌گان پرسش و پاسخ تصویری می‌پردازیم. سپس برجسته‌ترین روش‌های مطرح شده در هر دو رویکرد یادگیری عمیق و شبکه‌های از قبل آموزش دیده را بررسی خواهیم کرد.

#### ۳-۱ بررسی مجموعه داده‌گان مطرح این حوزه

در این بخش به معرفی مجموعه‌داده‌گان مشهور در حوزه پرسش و پاسخ تصویری می‌پردازیم و ویژگی‌های هر کدام را بررسی خواهیم کرد. در جدول ۳-۱ اطلاعات آماری این مجموعه‌داده‌گان به صورت خلاصه آمده‌است.

جدول ۳-۱: بررسی مجموعه دادگان در حوزه پرسش و پاسخ تصویری.

مجموعه داده	تعداد تصاویر	تعداد سوالات	سال انتشار
DAQUAR [۴۲]	۱۴۴۹	۱۲۴۶۸	۲۰۱۴
VQA v1 [۴]	۲۰۴۷۲۱	۶۱۴۱۶۳	۲۰۱۵
Visual Madlibs [۷۶]	۱۰۷۳۸	۳۶۰۰۰۱	۲۰۱۵
Visual7w [۸۰]	۴۷۳۰۰	۲۲۰۱۱۵۴	۲۰۱۶
VQA v2 [۱۹]	۲۰۴۷۲۱	۱۱۰۵۹۰۴	۲۰۱۷
CLEVR [۲۵]	۱۰۰۰۰۰	۸۵۳۵۵۴	۲۰۱۷
Tally-QA [۱]	۱۶۵۰۰۰	۳۰۶۹۰۷	۲۰۱۹
KVQA [۵۵]	۲۴۶۰۲	۱۸۳۰۰۷	۲۰۱۹

### ۳-۱-۱ دادگان DAQUAR

دادگان DAQUAR<sup>۱</sup> [۴۲] توسط مالدینوفسکی منتشر شده است. دادگان DAQUAR اولین مجموعه داده‌ای است که برای مسئله پرسش و پاسخ تصویری منتشر شده است. تصاویر از مجموعه داده NYU-Depth V2 [۵۹] گرفته شده است. اندازه این دادگان نسبتاً کوچک است و در مجموع ۱۴۴۹ تصویر دارد. دادگان DAQUAR شامل ۱۲۴۶۸ زوج پرسش و پاسخ با ۲۴۸۳ سوال منحصربه‌فرد است. برای تولید پرسش و پاسخ از دو روش مصنوعی و انسانی استفاده شده است. در روش مصنوعی پرسش و پاسخ‌ها به صورت خودکار از الگوهای موجود در جدول ۳-۲ تولید شده است. در روش دیگر از ۵ نفر خواسته شده تا پرسش و پاسخ تولید کنند. تعداد پرسش و پاسخ‌های آموزشی در این مجموعه داده ۶۷۹۴ و تعداد پرسش و پاسخ‌های تست ۵۶۴ است و به طور میانگین برای هر عکس تقریباً ۹ پرسش و پاسخ وجود دارد. این دادگان با مشکل بایاس<sup>۲</sup> روبه‌رو است زیرا تصاویر این مجموعه تنها مربوط به داخل خانه است و بیش از ۴۰۰ مورد وجود دارد که اشیایی مثل میز و صندلی در پاسخ‌ها تکرار شده است.

### ۳-۱-۲ دادگان VQA

دادگان Visual Question Answering v1 (VQA v1) [۴]<sup>۳</sup> یکی از پرکاربردترین دادگان در زمینه پرسش و پاسخ تصویری است. این دادگان شامل دو بخش است. یک بخش از تصاویر واقعی ساخته شده است که

<sup>۱</sup> Dataset for Question Answering on Real World Images

<sup>۲</sup> bias

<sup>۳</sup> <https://visualqa.org/>

جدول ۳-۲: الگوهای استفاده شده برای تولید سوال در دادگان DAQUAR. سوالات می‌تواند در مورد یک تصویر و یا مجموعه‌ای از تصاویر باشد [۴۲].

نمونه	الگو	توضیح	
How many cabinets are in image1?	How many {object} are in {image id}?	شماری	منفرد
How many gray cabinets are in image1?	How many {color} {object} are in {image id}?	شماری و رنگ	منفرد
Which type of the room is depicted in image1?	Which type of the room is depicted in {image id}?	نوع اتاق	منفرد
What is the largest object in image1?	What is the largest {object} in {image id}?	صفات عالی	منفرد
How many black bags?	How many {color} {object}?	شماری و رنگ	مجموعه‌ای
Which images do not have sofa?	Which images do not have {object}?	نفی نوع ۱	مجموعه‌ای
Which images are not bedroom?	Which images are not {room type}?	نفی نوع ۲	مجموعه‌ای
Which images have desk but do not have a lamp?	Which images have {object} but do not have a {object}?	نفی نوع ۳	مجموعه‌ای



شکل ۳-۱: چند نمونه از دادگان DAQUAR [۴۲]

VQA-real نام دارد و بخش دیگر با تصاویر کارتون ساخته شده است که با نام VQA-abstract از آن در مقالات یاد می‌شود.

بخش VQA-real به ترتیب شامل ۱۲۳۲۸۷ تصویر آموزشی و ۸۱۴۳۴ تصویر آزمایشی است که این تصاویر از دادگان MS-COCO [۳۵] تهیه شده است. برای جمع‌آوری پرسش و پاسخ از نیروی انسانی استفاده شده است. برای هر تصویر حداقل ۳ سوال منحصر به فرد وجود دارد و برای هر سوال ۱۰ پاسخ توسط کاربرهای منحصر به فرد جمع‌آوری شده است. این دادگان شامل ۶۱۴۱۶۳ سوال به صورت open-ended و چندگزینه‌ای است. در [۴] بررسی دقیقی در مورد نوع سوالات، طول سوالات و پاسخ‌ها انجام شده است.

بخش VQA-abstract به عنوان یک دادگان جداگانه و مکمل در کنار VQA-real قرار دارد. هدف از این دادگان از بین بردن نیاز به تجزیه و تحلیل تصاویر واقعی است تا مدل‌ها برای پاسخ به سوالات تمرکز خود را بر روی استدلال‌های سطح بالاتری بگذارند. تصاویر کارتون در این دادگان به صورت دستی توسط انسان‌ها و به وسیله‌ی رابط کاربری که از قبل آماده شده است؛ ساخته شده است. تصاویر می‌تواند دو حالت را نشان دهند: داخل خانه و خارج از خانه که هر کدام مجموعه متفاوتی از عناصر را شامل می‌شوند از جمله حیوانات،





Q: What shape is the bench seat ?

A: oval, semi circle, curved, curved, double curve, banana, curved, wavy, twisting, curved



Q: What color is the stripe on the train ?

A: white, white, white, white, white, white, white, white, white, white, white



Q: Where are the magazines in this picture ?

A: On stool, stool, on stool, on bar stool, on table, stool, on stool, on chair, on bar stool, stool

شکل ۳-۲: چند نمونه از دادگان VQA v1 - real [۴]



Q: Who looks happier ?.

A: old person, man, man, man, old man, man, man, man, man, grandpa



Q: Where are the flowers ?

A: near tree, tree, around tree, tree, by tree, around tree, around tree, grass, beneath tree, base of tree



Q: How many pillows ?

A: 1, 2, 2, 2, 2, 2, 2, 2, 2, 2

شکل ۳-۳: چند نمونه از دادگان VQA v1 - abstract [۴]

اشیا و انسان‌ها با حالت‌های مختلف. در مجموع ۵۰۰۰۰ تصویر ایجاد شده است. مشابه VQA-real، ۳ سوال برای هر تصویر (یعنی در کل ۱۵۰۰۰۰ سوال) و برای هر سوال ۱۰ پاسخ جمع‌آوری شده است. دادگان Visual Question Answering v2 (VQA v2) [۱۹] در سال ۲۰۱۷ پس از دادگان VQA v1 معرفی شد. دادگان VQA v2 نسبت به VQA v1 متوازن‌تر است و تعصبات زبانی در VQA v1 را کاهش داده است. اندازه‌ی دادگان VQA v2 تقریباً دو برابر دادگان VQA v1 است. در دادگان VQA v2 تقریباً برای هر سوال دو تصویر مشابه وجود دارد که پاسخ‌های متفاوتی برای سوال دارند.

### ۳-۱-۳ دادگان Visual Madlibs

دادگان Visual Madlibs [۷۶] شکل متفاوتی از پرسش و پاسخ را ارائه می‌دهد. برای هر تصویر جملاتی در نظر گرفته شده است و یک کلمه از آن که معمولاً مربوط به آدم، اشیا و فعالیت‌های نمایش داده شده در تصویر است؛ از جمله حذف شده و به جای آن جای خالی قرار گرفته است. پاسخ‌ها کلماتی هستند که این جملات را



شکل ۳-۴: چند نمونه از دادگان VQA v2 [۱۹]

تکمیل می‌کنند. برای مثال جمله ”دو [جای خالی] در پارک [جای خالی] بازی می‌کنند.“ در وصف یک تصویر بیان شده است که با دو کلمه ”مرد“ و ”فریزی“ می‌توان جاهای خالی را پر کرد. این دادگان شامل ۱۰۷۳۸ تصویر از دادگان MS-COCO [۳۵] و ۳۶۰۰۰۱ جمله با جای خالی است. جملات با جای خالی به طور خودکار و با استفاده از الگوهای از پیش تعیین شده تولید شده‌اند. پاسخ‌ها در این دادگان به هر دو شکل open-ended و چندگزینه‌ای است.

### ۳-۱-۴ دادگان Visual7w

دادگان Visual7W [۸۰] نیز بر اساس دادگان MS-COCO [۳۵] ساخته شده است. این دادگان شامل ۴۷۳۰۰ تصویر و ۳۲۷۹۳۹ جفت سوال و پاسخ است. این دادگان همچنین از ۱۳۱۱۷۵۶ پرسش و پاسخ چندگزینه‌ای تشکیل شده است که هر سوال ۴ گزینه دارد و تنها یکی از گزینه‌ها پاسخ صحیح سوال است. برای جمع‌آوری سوالات چندگزینه‌ای توسط انسان‌ها از پلتفرم آنلاین Amazon Mechanical Turk استفاده شده است. نکته‌ی حائز اهمیت در این دادگان این است که تمامی اشیایی که در متن پرسش یا پاسخ ذکر شده است، به نحوی به کادر محدودکننده‌ی آن شی در تصویر مرتبط شده است. مزیت این روش، رفع ابهام‌های موجود در متن است. همان‌طور که از نام این دادگان پیداست؛ سوالات آن با ۷ کلمه‌ی پرسشی که حرف اول آن w است شروع می‌شود. این ۷ کلمه شامل what ، where ، when ، who ، why ، how و which است. پرسش‌های Visual7W نسبت به به دادگان VQA v1 غنی‌تر و سخت‌تر است. همچنین پاسخ‌ها طولانی‌تر هستند.



1. This place is a park.
2. When I look at this picture, I feel competitive.
3. The most interesting aspect of this picture is the guys playing shirtless.
4. One or two seconds before this picture was taken, the person caught the frisbee.
5. One or two seconds after this picture was taken, the guy will throw the frisbee.
6. Person A is wearing blue shorts.
7. Person A is in front of person B.
8. Person A is blocking person B.
9. Person B is a young man wearing an orange hat.
10. Person B is on a grassy field.
11. Person B is holding a frisbee.
12. The frisbee is white and round.
13. The frisbee is in the hand of the man with the orange cap.
14. People could throw the frisbee.
15. The people are playing with the frisbee.

شکل ۳-۵: یک نمونه از دادگان Visual Madlibs [۷۶]



Q: What endangered animal is featured on the truck?

A: A bald eagle.  
A: A sparrow.  
A: A humming bird.  
A: A raven.



Q: Where will the driver go if turning right?

A: Onto 24 1/2 Rd.  
A: Onto 25 1/2 Rd.  
A: Onto 23 1/2 Rd.  
A: Onto Main Street.



Q: When was the picture taken?

A: During a wedding.  
A: During a bar mitzvah.  
A: During a funeral.  
A: During a Sunday church service.



Q: Who is under the umbrella?

A: Two women.  
A: A child.  
A: An old man.  
A: A husband and a wife.



Q: Why was the hand of the woman over the left shoulder of the man?

A: They were together and engaging in affection.  
A: The woman was trying to get the man's attention.  
A: The woman was trying to scare the man.  
A: The woman was holding on to the man for balance.



Q: How many magnets are on the bottom of the fridge?

A: 5.  
A: 2.  
A: 3.  
A: 4.



Q: Which pillow is farther from the window?



Q: Which step leads to the tub?



Q: Which is the small computer in the corner?



Q: Which item is used to cut items?

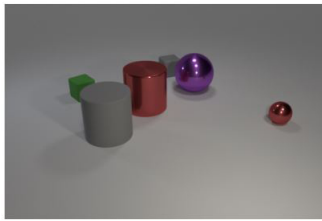


Q: Which doughnut has multicolored sprinkles?

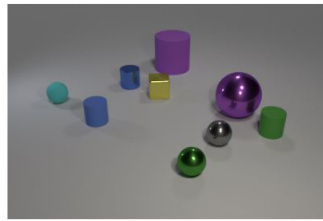


Q: Which man is wearing the red tie?

شکل ۳-۶: چند نمونه از دادگان Visual7W [۸۰]. ردیف اول، پاسخ‌های سبز، پاسخ صحیح هستند و پاسخ‌های قرمز پاسخ‌های نادرست تولید شده توسط انسان است. ردیف دوم، کادر زرد جواب صحیح است و کادرهای قرمز پاسخ‌های اشتباه انسانی است.

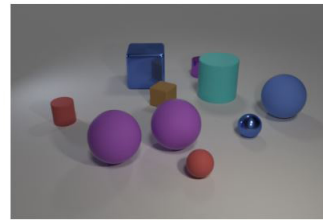


Q: How big is the gray rubber object that is behind the big shiny thing behind the big metallic thing that is on the left side of the purple ball?  
A: small



Q: There is a tiny rubber thing that is the same color as the metal cylinder; what shape is it?  
A: cylinder

Q: What is the shape of the tiny green thing that is made of the same material as the large cylinder?  
A: cylinder



Q: There is a small ball that is made of red rubber sphere the same material as the large block; what color is it?  
A: blue

Q: Is the size of the ball that is made of red rubber sphere the same as the purple metal thing?  
A: yes

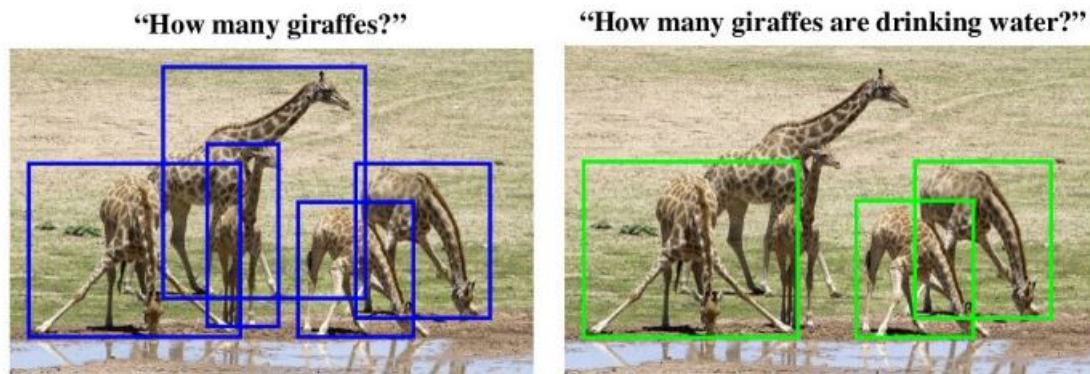
شکل ۳-۷: چند نمونه از دادگان CLEVR [۲۵].

### ۳-۱-۵ دادگان CLEVR

دادگان CLEVR [۲۵] یک دادگان برای ارزیابی درک بصری سیستم‌های پرسش و پاسخ تصویری است. تصاویر این دادگان با استفاده از سه شی استوانه، کره و مکعب تولید شده است. برای هر کدام از این اشیا دو اندازه متفاوت، دو جنس متفاوت و هشت رنگ مختلف در نظر گرفته شده است. سوالات هم به طور مصنوعی بر اساس مکانی که اشیا در تصویر قرار گرفته‌اند؛ ایجاد شده‌است. سوالات در CLEVR به گونه‌ای طراحی شده است که جنبه‌های مختلف استدلال بصری توسط سیستم‌های پرسش و پاسخ تصویری را مورد ارزیابی قرار می‌دهد از جمله شناسایی ویژگی، شمارش اشیا، مقایسه، روابط مکانی اشیا و عملیات منطقی. در این دادگان مکان تصاویر نیز با استفاده از یک مستطیل مشخص شده است.

### ۳-۱-۶ دادگان Tally-QA

در سال ۲۰۱۹، دادگان Tally-QA [۱] منتشر شد که بزرگ‌ترین دادگان پرسش و پاسخ تصویری برای شمارش اشیا است. اکثر مجموعه داده‌های شمارش اشیا در پرسش و پاسخ تصویری دارای سوالات ساده هستند که برای پاسخ دادن به این سوال‌ها تنها کافی است که اشیا در تصویر تشخیص داده شوند. بنابراین، این موضوع باعث ایجاد مجموعه داده‌ی Tally-QA شد که علاوه بر سوالات ساده، سوالات پیچیده را نیز در بر می‌گیرد که برای پاسخ دادن به آن‌ها به استدلال بیشتری از تشخیص اشیا نیاز است. تعداد سوالات ساده در Tally-QA برابر با ۲۱۱۴۳۰ و تعداد سوالات پیچیده برابر با ۷۶۴۷۷ است. سوالات ساده این دادگان از مجموعه داده‌های



شکل ۳-۸: چند نمونه از دادگان Tally-QA [۱]. عکس سمت چپ یک نمونه از سوالات ساده و عکس سمت راست یک نمونه از سوالات پیچیده است.

دیگری (VQA v2 [۱۹] و Visual Genome [۳۱]) برداشته شده است و سوالات پیچیده با استفاده از ۸۰۰ کاربر انسانی از طریق پلتفرم آنلاین Amazon Mechanical Turk جمع‌آوری شده‌است. دادگان Tally-QA به سه بخش آموزش و تست - ساده و تست - پیچیده تقسیم می‌شود. بخش تست - ساده تنها شامل سوالات ساده و بخش تست - پیچیده تنها دارای سوالات پیچیده‌ای است که از Amazon Mechanical Turk جمع‌آوری شده است.

### ۳-۱-۷ دادگان KVQA

دادگان KVQA [۵۵] که مخفف Knowledge-based Visual Question Answering است در سال ۲۰۱۹ طراحی شده است به طوری که بر خلاف مجموعه داده‌های قبلی، برای پیدا کردن پاسخ سوالات نیاز به دانش خارجی دارد. بدین منظور این دادگان شامل ۱۸۳ هزار پرسش و پاسخ در مورد ۱۸ هزار شخص معروف شامل ورزشکاران، سیاستمداران و هنرمندان است. اطلاعات و تصاویر مرتبط با این اشخاص از Wikidata و Wikipedia استخراج شده است. دادگان KVQA شامل ۲۴ هزار تصویر است. این دادگان به صورت تصادفی به سه بخش آموزش، ارزیابی و آزمون به ترتیب با نسبت‌های ۰.۷، ۰.۲ و ۰.۱ تقسیم شده است. تنوع پرسش و پاسخ‌ها در KVQA به گونه‌ای در نظر گرفته شده است که مشکل همیشگی بایاس در مجموعه داده‌های پرسش و پاسخ تصویری، در این مجموعه داده وجود نداشته باشد.





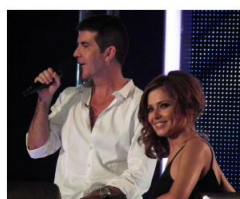
(a) *Wikipedia caption*: Khan with United States Secretary of State Hillary Clinton in 2009.

Q: Who is to the left of Hillary Clinton? (*spatial*)

A: **Aamir Khan**

Q: Do all the people in the image have a common occupation? (*multi-entity, intersection, 1-hop, Boolean*)

A: **No**



(b) *Wikipedia caption*: Cheryl alongside Simon Cowell on The X Factor, London, June 2010.

Q: What is the age gap between the two people in the image? (*multi-entity, subtraction, 1-hop*)

A: **24 years**

Q: How many people in this image were born in United Kingdom? (*1-hop, multi-entity, counting*)

A: **2**



(c) *Wikipedia caption*: BRICS leaders at the G-20 summit in Brisbane, Australia, 15 November 2014

Q: Were all the people in the image born in the same country? (*Boolean, multi-entity, intersection*)

A: **No**

Q: Who is the founder of the political party to which person second from left belongs to? (*spatial, multi-hop*)

A: **Syama Prasad Mookerjee**



(d) *Wikipedia caption*: Serena Williams and Venus Williams, Australian Open 2009.

Q: Who among the people in the image is the eldest? (*multi-entity, comparison*)

A: **Person in the left**

Q: Who among the people in the image were born after the end of World War II? (*multi-entity, multi-relation, comparison*)

A: **Both**

شکل ۳-۹: چند نمونه از دادگان KVQA [۵۵]

## ۲-۳ تقویت دادگان در مسئله پرسش و پاسخ تصویری

با توسعه سریع شبکه‌های عصبی عمیق مسئله پرسش و پاسخ تصویری به موفقیت‌های بزرگی دست یافته است. مطالعات نشان می‌دهد که عملکرد شبکه‌های عصبی عمیق به میزان داده‌های آموزشی بستگی دارد و همیشه از داده‌های آموزشی بیشتر سود می‌برند. یکی از ترفندهای اصلی در شبکه‌های عصبی عمیق تقویت داده<sup>۴</sup> است که به طور گسترده در بسیاری از مسائل پردازش تصویر و بینایی ماشین مورد استفاده قرار می‌گیرد. اما مقالات کمی وجود دارد که مسئله تقویت داده را در پرسش و پاسخ تصویری بررسی کرده‌اند. یکی از چالش‌های تقویت داده در مسئله پرسش و پاسخ تصویری این است که هیچ یک از روش‌های تقویت داده مبتنی بر تصویر مانند چرخش<sup>۵</sup> نمی‌توانند مستقیماً بر روی مسئله پرسش و پاسخ تصویری اعمال شود زیرا ساختار معنایی آن حفظ نخواهد شد. به عنوان مثال با چرخش یک تصویر ممکن است پرسش و پاسخ مرتبط با آن (مانند «ماشین در سمت چپ یا راست سطل زباله است؟») دیگر درست نباشد.

در [۲۸] برای اولین بار دو روش برای تقویت داده در مسئله پرسش و پاسخ تصویری پیشنهاد شد. در روش اول برای تولید پرسش و پاسخ از الگو استفاده می‌شود. برای تولید الگو از حاشیه‌نویسی<sup>۶</sup> موجود

<sup>۴</sup> data augmentation  
<sup>۵</sup> rotation  
<sup>۶</sup> annotation

در دادگان استفاده می‌شود. با استفاده از این روش ۴ نوع سوال تولید می‌شود: (۱) سوالات بله و خیر (۲) سوالات شمارشی (۳) سوالات درباره تشخیص شی، صحنه و یا فعالیت (۴) سوالات درباره تشخیص ورزش. برای مثال برای تولید سوالات بله و خیر، با استفاده از حاشیه‌نویسی موجود در دادگان لیستی از اشیا موجود در تصویر آماده می‌شود. سپس اگر محدوده مربوط به اشیا بزرگتر از ۲۰۰۰ پیکسل باشد، سوالی مانند «آیا [شی] در تصویر وجود دارد؟» تولید می‌شود که پاسخ آن هم «بله» است. به همین ترتیب با استفاده از دانشی که از دادگان می‌توان بدست آورد؛ برای سایر انواع سوالات الگویی برای تولید سوال و پاسخ آن تولید می‌شود. یکی از مشکلات این روش برای تقویت داده این است که سوالات تولید شده انعطاف‌پذیر نیستند و ممکن است شباهت چندانی به سوالات موجود در دادگان نداشته باشند. به همین علت، روش دیگری در [۲۸] مبتنی بر LSTM برای تولید سوال برای هر تصویر پیشنهاد شده است. این شبکه از دو لایه LSTM تشکیل شده است که هر کدام دارای ۱۰۰۰ واحد مخفی است و پس از آن‌ها نیز دو لایه کاملاً متصل که هر کدام ۷۰۰۰ نورون مخفی دارند (برابر با تعداد واژگان) ساخته شده است. برای تولید سوال، در ابتدا توکن<sup>۷</sup> شروع سوال به همراه ویژگی‌های تصویر به شبکه داده می‌شود. برای هر تصویر ۳۰ سوال تولید می‌شود که تنها سه تا از پرتکرارترین سوالات نگه داشته می‌شود. برای پیدا کردن جواب سوال‌های تولید شده توسط شبکه LSTM از یک شبکه‌ی ساده MLP<sup>۸</sup> که در [۲۶] پیشنهاد شده است؛ استفاده شده است. در [۲۸] نشان دادند که استفاده از این دو روش برای تقویت داده‌ها منجر به بهبود عملکرد روش‌های موجود برای حل مسئله پرسش و پاسخ تصویری می‌شود.

اخیراً در [۶۶] برای تقویت داده روشی مبتنی بر تولید نمونه‌های خصمانه<sup>۹</sup> پیشنهاد شده است که بر خلاف کارهای قبلی، تقویت داده هم برای تصاویر و هم برای سوالات انجام می‌شود.

### ۳-۳ رویکرد یادگیری عمیق

اکثر روش‌های پیشنهاد شده در رویکرد یادگیری عمیق دارای سه فاز هستند. فاز اول این فرآیند استخراج ویژگی از تصویر و سوالات است که راه‌حل‌های موفق در این فاز ریشه در روزهای باشکوه یادگیری عمیق دارد زیرا بیشتر راه‌حل‌های موفق در این حوزه از مدل‌های یادگیری عمیق استفاده می‌کنند مانند CNN ها برای استخراج

token<sup>۷</sup>  
Multi Layer Perceptron<sup>۸</sup>  
adversarial examples<sup>۹</sup>

ویژگی از تصویر و RNN ها و انواع آن (LSTM و GRU) برای استخراج ویژگی از سوالات. در فاز دوم که مهم‌ترین و اصلی‌ترین فاز می‌باشد، ویژگی‌های استخراج شده از تصویر و سوال باهم ترکیب می‌شوند. سپس از ترکیب ویژگی‌ها برای پیش‌بینی پاسخ نهایی در فاز سوم استفاده می‌شود. در ادامه این بخش جزئیات هر فاز را بررسی می‌کنیم.

### ۳-۳-۱ فاز ۱: استخراج ویژگی از تصویر و سوال

استخراج ویژگی از تصویر و سوال مرحله‌ی مقدماتی در پرسش و پاسخ تصویری است. ویژگی تصویر، تصویر را به عنوان یک بردار عددی توصیف می‌کند تا بتوان به راحتی عملیات‌های مختلف ریاضی را بر روی آن اعمال کرد. روش‌های زیادی وجود دارد که به صورت مستقیم از تصویر ویژگی استخراج می‌کنند مانند SIFT، تبدیل HAAR و HOG. اما با ظهور شبکه‌های یادگیری عمیق، نیاز به استخراج ویژگی به صورت مستقیم از بین رفت زیرا این شبکه‌ها قادر به یادگیری ویژگی هستند. آموزش مدل‌های یادگیری عمیق به منابع محاسباتی گران قیمت و مجموعه‌داده‌های بزرگ نیاز دارد. از این رو، استفاده از مدل‌های شبکه عصبی عمیق از قبل آموزش دیده، استخراج ویژگی از تصاویر را به راحتی امکان‌پذیر می‌کند.

یکی از بهترین شبکه‌های عصبی برای استخراج ویژگی از تصویر، شبکه‌های عصبی پیچشی هستند. در جدول ۱-۲ چند نمونه از برجسته‌ترین شبکه‌های عصبی پیچشی که بر روی دادگان ImageNet [۱۳] آموزش داده شده‌اند؛ آورده شده است. بیشتر مدل‌های ارائه‌شده در پرسش و پاسخ تصویری از این شبکه‌های عصبی پیچشی استفاده می‌کنند تا محتوای تصویری خود را به بردارهایی عددی تبدیل کنند. جدول ۳-۳ لیستی از مدل‌های استفاده شده برای حل مسئله پرسش و پاسخ تصویری را نشان می‌دهد و مشخص می‌کند که هر کدام از این مدل‌ها برای استخراج ویژگی از تصویر از کدام یک از شبکه‌های عصبی پیچشی موجود در جدول ۱-۲ بهره می‌برد. همان‌طور که واضح است VGGNet و ResNet به طور گسترده‌ای در سیستم‌های پرسش و پاسخ تصویری مورد استفاده قرار گرفته‌اند. یکی از دلایلی که محققان VGGNet را ترجیح می‌دهند این است که ویژگی‌هایی را استخراج می‌کند که عمومیت بیشتری دارد و برای مجموعه‌داده‌هایی غیر از ImageNet که این مدل‌ها بر روی آن‌ها آموزش داده می‌شوند، موثرتر هستند. دلایل دیگر شامل همگرایی سریع در fine-tuning و پیاده‌سازی ساده در مقایسه با GoogLeNet و ResNet است. نکته‌ی قابل توجه دیگر در جدول ۳-۳ روند مهاجرت از VGGNet به ResNet در مقالات اخیر است. زیرا در سال‌های اخیر، منابع محاسباتی کافی با هزینه مناسب در دسترس محققان می‌باشد.



جدول ۳-۳: شبکه‌های عصبی پیچشی استفاده شده در مدل‌های پرسش و پاسخ تصویری.

مدل پرسش و پاسخ تصویری	AlexNet	VGGNet	GoogleNet	ResNet
Image_QA[۵۲]		✓		
Talk_to_Machine[۱۷]			✓	
VQA[۴]		✓		
Vis_Madlibs[۷۶]	✓	✓		
VIS + LSTM[۵۱]		✓		
Ahab[۷۰]		✓		
ABC-CNN[۹]		✓		
Comp_QA[۳]		✓		
DPPNet[۴۶]		✓		
Answer_CNN[۴۰]		✓		
VQA-Caption[۳۶]		✓		
Re_Baseline[۲۴]				✓
MCB[۱۶]				✓
SMem-VQA[۷۴]			✓	
Region_VQA[۵۷]		✓		
Vis7W[۸۰]		✓		
Ask_Neuron[۴۳]	✓	✓	✓	✓
SCMC[۸]				✓
HAN[۴۱]				✓
StrSem[۷۸]		✓		
AVQAN[۵۴]				✓
CMF[۳۲]				✓
EnsAtt[۳۷]				✓
MetaVQA[۶۷]				✓
DA-NTN[۵]				✓
QGHC[۸]				✓
QTA[۵۶]				✓
WRAN[۴۸]				✓
QAR[۶۹]				✓

جدول ۳-۴: تعبیه کلمات استفاده شده در مدل‌های پرسش و پاسخ تصویری.

مدل پرسش و پاسخ تصویری	one-hot	CBOW	Skip-gram/Word2vec	GloVe	CNN	LSTM	GRU
Image_QA[۵۲]			✓				
Talk_to_Machine[۱۷]						✓	
VQA[۴]		✓					
Vis_Madlibs[۷۶]			✓				
VIS + LSTM[۵۱]						✓	
ABC-CNN[۹]						✓	
Comp_QA[۳]						✓	
DPPNet[۴۶]							✓
Answer_CNN[۴۰]					✓		
VQA-Caption[۳۶]						✓	
Re_Baseline[۲۴]			✓				
MCB[۱۶]						✓	
SMem-VQA[۷۴]		✓					
Region_VQA[۵۷]			✓				
Vis7W[۸۰]	✓						
Ask_Neuron[۴۳]		✓			✓	✓	✓
SCMC[۸]					✓		
HAN[۴۱]						✓	
StrSem[۷۸]						✓	
AVQAN[۵۴]	✓						
CMF[۳۲]				✓		✓	
EnsAtt[۳۷]				✓			
MetaVQA[۶۷]				✓			✓
DA-NTN[۵]							✓
QGHG[۸]							✓
WRAN[۴۸]							✓
QAR[۶۹]				✓			

مدل‌های مختلف در مسئله پرسش و پاسخ تصویری از تعبیه کلمات متفاوتی برای تولید بردار ویژگی سوال‌ها استفاده کرده‌اند. جدول ۳-۴ لیستی از مدل‌های پرسش و پاسخ تصویری به همراه تعبیه کلمات استفاده شده در آن‌ها را نمایش می‌دهد. با بررسی جدول ۳-۴ مشاهده می‌کنیم که محققان حوزه‌ی پرسش و پاسخ تصویری ترجیح می‌دهند؛ برای استخراج ویژگی از متن و بازنمایی آن از LSTM استفاده کنند. آن‌ها معتقد هستند که RNN ها عملکرد بهتری نسبت به روش‌های مستقل از دنباله‌ی کلمات مانند word2vec دارند. اما آموزش RNN ها نیاز به داده‌های برچسب خورده‌ی زیادی دارد.

## ۳-۳-۲ فاز ۲: بازنمایی مشترک تصویر و سوال

در گام اول پرسش و پاسخ تصویری، تصویر و سوال به طور مستقل پردازش می‌شوند تا از آن‌ها ویژگی استخراج شود. روش‌های مختلف برای انجام این کار، در بخش ۳-۳-۱ به تفصیل بررسی شد. در گام بعدی، این ویژگی‌ها باید به یک فضای مشترک ترسیم شوند و یا به عبارتی ترکیب شوند تا آماده گام آخر (تولید پاسخ) شوند. در ادامه این بخش، به مرور روش‌های ترکیب ویژگی‌های استخراج شده از سوال و تصویر می‌پردازیم.

## ۳-۳-۲-۱ روش‌های پایه

ساده‌ترین و پایه‌ای‌ترین روش‌ها برای ترکیب ویژگی‌ها concatenation، جمع متناظر ویژگی‌ها<sup>۱۰</sup> و ضرب متناظر ویژگی‌ها<sup>۱۱</sup> است. مالینوفسکی در [۴۳] این سه روش را امتحان کرده است و دریافت کرد که ضرب متناظر ویژگی‌ها منجر به دقت بالاتری می‌شود. یافته مهم دیگر مالینوفسکی این است که نرمال‌سازی L2 ویژگی‌های تصویر، تأثیر قابل توجهی دارد به خصوص در روش‌های concatenation و جمع متناظر ویژگی‌ها. با توجه به نتایج آن‌ها، جمع متناظر ویژگی‌ها پس از نرمال‌سازی از دقت بالاتری برخوردار است.

روش کلاسیک دیگر برای یافتن رابطه بین دو بردار که ریشه آن در علم آمار است، روش CCA<sup>۱۲</sup> است که برای ترکیب ویژگی‌های تصویر و سوال در VQA استفاده شده است. CCA بازنمایی مشترک بین بردار تصویر و بردار سوال را پیدا می‌کند. CCA یک نسخه نرمالیزه شده به نام nCCA<sup>۱۳</sup> نیز دارد که توسط [۱۸] پیشنهاد شده است. در [۷۶] و [۶۸] از هر دو مدل CCA و nCCA برای ترکیب بردارهای ویژگی سوال و تصویر استفاده کردند و دریافتند که روش nCCA به ویژه در مورد سوالات چندگزینه‌ای عملکرد بهتری دارد.

## ۳-۳-۲-۲ روش‌های مبتنی بر شبکه‌های عصبی

در این روش‌ها، محققان شبکه‌های عصبی را با لایه‌های خاص برای ترکیب ویژگی‌های تصویر و سوال آموزش می‌دهند. ساختار و عملکرد این لایه ممکن است برای مدل‌های مختلف پیشنهاد شده متفاوت باشد. از این رو، مدل‌های پیشنهاد شده با این روش برای مسئله پرسش و پاسخ تصویری بسیار زیاد و متفاوت است. بنابراین در ادامه چند نمونه از شبکه‌های پیشنهاد شده را معرفی می‌کنیم.

<sup>۱۰</sup> element-wise addition<sup>۱۱</sup> element-wise multiplication<sup>۱۲</sup> Analysis Correlation Canonical<sup>۱۳</sup> Analysis Correlation Canonical normalized

در [۱۷] برای ترکیب ویژگی‌های تصویر و سوال از یک لایه استفاده شده است که ساختار اصلی آن تابع فعالساز غیرخطی  $\tanh$  است. پس از جمع متناظر ویژگی‌های تصویر و سوال با هم، حاصل به این لایه داده می‌شود تا ویژگی‌ها با هم ترکیب شوند. تابع اعمال شده در این لایه در عبارت ۱-۳ آورده شده است.

$$g(x) = 1/7159 \tanh\left(\frac{2}{3}x\right) \quad (1-3)$$

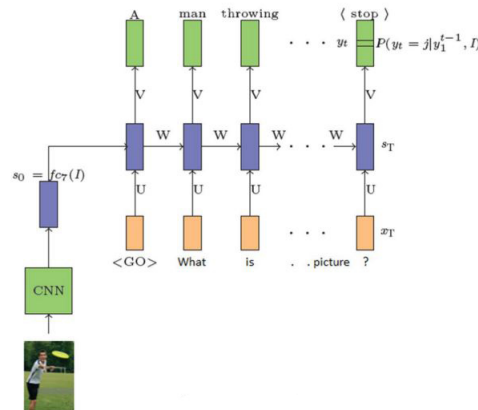
در [۴۰] علاوه بر این که برای استخراج ویژگی از تصویر و سوال از CNN استفاده شده است؛ برای ترکیب ویژگی‌ها نیز از شبکه عصبی پیچشی استفاده شده است که آن را multimodel CNN نامیده‌اند. برای انجام کانولوشن در multimodel CNN در هر پنجره بازنمایی کل تصویر به همراه بازنمایی دو کلمه متوالی از سوال در نظر گرفته می‌شود.

نویسندگان مقاله [۴۶] معتقدند که ثابت بودن پارامترهای شبکه عصبی به اندازه کافی برای مسئله پرسش و پاسخ تصویری قدرتمند نیستند. به همین دلیل آن‌ها پس از شبکه VGGNet سه لایه کاملاً متصل قرار می‌دهند که پارامترهای دومین لایه کاملاً متصل را متغیر و متناسب با سوال ورودی تنظیم می‌کنند. از این رو یک لایه به نام DPPN<sup>۱۴</sup> طراحی کردند که از یک GRU برای بازنمایی سوال استفاده می‌کند و سپس با استفاده از یک تابع هش، پارامترهای لایه دوم کاملاً متصل را محاسبه می‌کند.

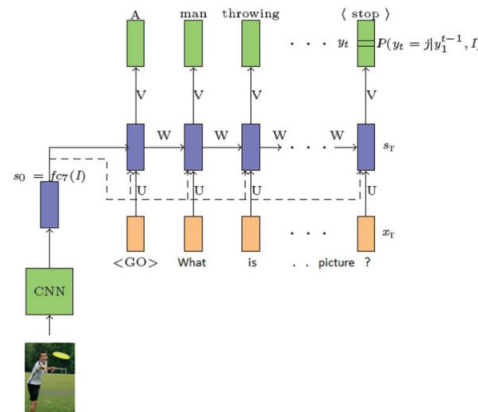
برای ترکیب ویژگی‌های تصویر و سوال از ایده‌ی یادگیری باقی‌مانده<sup>۱۵</sup> (که در شبکه ResNet به کار گرفته شد) در مسئله پرسش و پاسخ تصویری نیز استفاده می‌شود. برای مثال در [۲۹] شبکه MRN<sup>۱۶</sup> بر مبنای همین ایده برای ترکیب ویژگی‌های تصویر و سوال در مسئله پرسش و پاسخ تصویری پیشنهاد شده است.

یکی از جدیدترین روش‌های ترکیب ویژگی در مسئله پرسش و پاسخ تصویری با استفاده از شبکه‌های عصبی عمیق، معماری رمزگذار-رمزگشا<sup>۱۷</sup> است. در این روش، سوال رمزگذاری شده به همراه تصویر رمزگذاری شده به رمزگشا داده می‌شود (معمولاً از LSTM به عنوان رمزگشا استفاده می‌شود). سپس برای تولید پاسخ صحیح آموزش داده می‌شود. معماری این روش به صورت کلی دو حالت می‌تواند داشته باشد. در حالت اول، بازنمایی تصویر به عنوان اولین کلمه از دنباله کلمات سوال به کدگشا داده می‌شود (شکل ۳-۱۰). مدل‌های پیشنهاد شده در [۵۲] و [۸۰] از این نوع هستند. در حالت دوم، بازنمایی تصویر در هر گام زمانی به

<sup>۱۴</sup>Dynamic Parameter Prediction Network<sup>۱۵</sup>residual learning<sup>۱۶</sup>Multimodal Residual Network<sup>۱۷</sup>encoder-decoder architecture



شکل ۳-۱۰: حالت اول معماری رمزگذار-رمزگشا در پرسش و پاسخ تصویری



شکل ۳-۱۱: حالت دوم معماری رمزگذار-رمزگشا در پرسش و پاسخ تصویری

LSTM داده می‌شود (شکل ۳-۱۱). در [۴۳] از این روش استفاده شده است. در [۵۴] علاوه بر رمزگذاری تصویر و رمزگذاری سوال که به عنوان ورودی به رمزگشا داده می‌شود، رمزگذاری دیگری به نام question mood به عنوان ورودی سوم به رمزگشا داده می‌شود تا یک صفت احساسی همراه با پاسخ تولید شود. در مقاله [۷۱] روشی که پیشنهاد شده است مطابق با شکل ۳-۱۰ است با این تفاوت که ورودی اولین گام زمانی LSTM شامل بردار ویژگی‌های تصویر، تعبیه جمله توصیف‌کننده تصویر و یک بردار از دانشی که با توجه به سوال از منابع خارجی استخراج شده است، می‌باشد. این روش برای پاسخ دادن به سوالات «چرا» بسیار مناسب است.

## ۳-۲-۳-۳ روش‌های مبتنی بر مکانیزم توجه

در ۵ سال گذشته، روش‌های بسیاری در مسئله پرسش و پاسخ تصویری مطرح شده است که اساس کار آن‌ها بر پایه مکانیزم توجه<sup>۱۸</sup> است. مدل‌های مبتنی بر مکانیزم توجه به ناحیه‌هایی از تصاویر که مربوط به سوال است، توجه می‌کنند. مدل‌های موجود در این رویکرد یا به تصویر و یا به سوال و یا به هر دو توجه می‌کنند. در ادامه این بخش چند نمونه از برجسته‌ترین روش‌های پیشنهادشده بر پایه مکانیزم توجه در مسئله پرسش و پاسخ تصویری را بررسی می‌کنیم.

در [۷۵] مدلی به نام Stacked Attention Network (SAN) پیشنهاد شده که ایده‌ی اصلی آن این است که ابتدا از سوال، یک بازنمایی معنایی و مفهومی استخراج می‌شود. سپس از آن به عنوان یک کوئری برای پیدا کردن مناطقی از تصویر که مرتبط با سوال است؛ استفاده می‌شود. غالباً در مسئله پرسش و پاسخ تصویری نیاز است تا چندین مرحله استدلال صورت بگیرد. بنابراین در این شبکه از چندین لایه برای جستجو در تصویر استفاده می‌شود تا به تدریج به جواب مورد نظر برسد.

روش پیشنهاد شده در [۳۹] همزمان هم به تصویر و هم به سوال توجه می‌کند. این روش دارای دو ویژگی مهم است. ویژگی اول بازنمایی سلسله‌مراتبی سوال و ویژگی دوم مکانیزم توجه همزمان<sup>۱۹</sup> می‌باشد. روند کلی در مکانیزم توجه همزمان به این صورت است که از بازنمایی تصویر برای محاسبه توجه سوال استفاده می‌شود و به طور متقابل از بازنمایی سوال برای محاسبه توجه تصویر استفاده می‌شود. بنابراین در [۳۹] ابتدا برای سوال یک بازنمایی سلسله‌مراتبی محاسبه می‌شود که شامل تعبیه کلمات، تعبیه عبارات و تعبیه جمله است. سپس مکانیزم توجه همزمان در هر کدام از این سه سطح هم برای سوال و هم برای تصویر انجام می‌شود و پاسخ نهایی بر اساس خروجی‌های حاصل از این مرحله بدست می‌آید.

یکی از نوآوری‌های اخیر در مکانیزم توجه، توجه سخت<sup>۲۰</sup> است که مالینوفسکی در [۴۱] برای حل مسئله پرسش و پاسخ تصویری استفاده کرده است. توجه نرم<sup>۲۱</sup> که عموماً از لفظ توجه برای آن استفاده می‌شود، با محاسبه میانگین وزندار مشخص می‌کند که به کدام مناطق از ویژگی‌ها توجه بیشتری شود و به کدام بخش‌ها توجه کمتری شود. اما در توجه سخت، از ویژگی‌ها نمونه‌برداری<sup>۲۲</sup> می‌شود و یک یا چند ویژگی در خروجی ظاهر می‌شود. البته این نمونه‌برداری براساس یک توضیح احتمالاتی انجام می‌شود که ویژگی‌های معنادار

<sup>۱۸</sup> attention mechanism  
<sup>۱۹</sup> coattention mechanism  
<sup>۲۰</sup> hard attention  
<sup>۲۱</sup> soft attention  
<sup>۲۲</sup> sampling

جدول ۳-۵: بررسی رویکرد پیش‌بینی پاسخ در چند نمونه از مدل‌های پرسش و پاسخ تصویری.

تولید	طبقه‌بندی	مدل پرسش و پاسخ تصویری
✓		Talk_to_Machine[۱۷]
✓	✓	VQA[۴]
	✓	HieCoAttention[۳۹]
	✓	MCB[۱۶]
✓	✓	Ask_Neuron[۴۳]
	✓	Mutan[۶]
	✓	MCAN[۷۷]
	✓	AnswerAll[۵۸]

احتمال بیشتری دارند که در نمونه‌برداری انتخاب شوند. بنابراین در توجه سخت، با نمونه‌برداری اطلاعات ناخواسته حذف می‌شوند. از طرفی حالت قطعی بودن برخلاف توجه نرم در توجه سخت از بین می‌رود و این باعث می‌شود که این فرایند مشتق ناپذیر باشد و برای آموزش این مدل‌ها نتوان از روش کاهش گرادیان برای بهینه‌سازی مدل استفاده کرد. مالینوفسکی [۴۱] از ایده توجه سخت برای حذف المان‌هایی که اهمیت کمتری در ترکیب ویژگی‌های تصویر و سوال دارند، استفاده کرده است.

### ۳-۳-۳ فاز ۳: پیش‌بینی پاسخ

در این فاز برای بدست آوردن پاسخ، به طور کلی از دو رویکرد طبقه‌بندی<sup>۲۳</sup> و تولید<sup>۲۴</sup> استفاده می‌شود. در رویکرد طبقه‌بندی مجموعه‌ای از پیش تعیین شده از پاسخ‌های کاندید آماده می‌شود و هر کدام از پاسخ‌های کاندید به عنوان یک کلاس در نظر گرفته می‌شود. بنابراین در مدل‌های پیشنهادی برای مسئله پرسش و پاسخ تصویری که از رویکرد طبقه‌بندی استفاده می‌کنند، در آخرین لایه از یک تابع softmax استفاده می‌شود و پاسخی که بیشترین احتمال را داشته باشد به عنوان پاسخ پیش‌بینی شده مدل در نظر گرفته می‌شود. در رویکرد طبقه‌بندی برای بدست آوردن مجموعه‌ای از پاسخ‌های کاندید، معمولاً  $n$  پاسخی که بیشترین تکرار را در دادگان داشته‌اند را در نظر می‌گیرند. در رویکرد تولید پاسخ، معمولاً از بازنمایی مشترک تصویر و سوال استفاده می‌شود و به کمک LSTM یک جمله به عنوان پاسخ در خروجی تولید می‌شود.

<sup>۲۳</sup> classification  
<sup>۲۴</sup> generation

فصل ۳. مروری بر کارهای مرتبط ۳-۴. رویکرد مدل‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر

در جدول ۳-۵ رویکرد پیش‌بینی پاسخ استفاده شده در چند نمونه از مدل‌های پرسش و پاسخ تصویری آورده شده است. همان‌طور که واضح است بیشتر مدل‌ها از رویکرد طبقه‌بندی برای پیش‌بینی پاسخ استفاده کرده‌اند.

## ۳-۴ رویکرد مدل‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر

در سال‌های اخیر شاهد ظهور شبکه‌های از قبل آموزش دیده تنها بر روی داده‌های تصویری مثل ResNet [۲۱] و یا تنها بر روی داده‌های متنی مانند BERT [۱۵]، GPT-2 [۵۰] و GPT-3 [۷] بوده‌ایم. استفاده از این شبکه‌ها منجر به بهبود مسائل موجود در بینایی ماشین و پردازش زبان‌های طبیعی شده است. با الهام از این موضوع، مدل‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر<sup>۲۵</sup> نیز ایجاد شدند که هدف آن‌ها بازنمایی مشترک داده‌های تصویری و داده‌های زبانی است. بنابراین می‌توان از این شبکه‌ها برای بهبود عملکرد مسائل مشترک بین بینایی ماشین و پردازش زبان‌های طبیعی مانند پرسش و پاسخ تصویری نیز استفاده کرد. معماری شبکه‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر به طور کلی به دو دسته تک جریانی<sup>۲۶</sup> و دو جریانی<sup>۲۷</sup> تقسیم می‌شود. در ادامه به بحث و بررسی هر یک از این دسته‌ها می‌پردازیم.

### ۳-۴-۱ معماری تک جریانی

پایه و اساس این معماری شبیه معماری مدل BERT [۱۵] است که رمزگذاری متن<sup>۲۸</sup> و رمزگذاری تصویر<sup>۲۹</sup> را به طور همزمان انجام می‌دهد. در واقع برای یادگیری بازنمایی متن و تصویر از یک رمزگذار<sup>۳۰</sup> استفاده می‌کند. بنابراین ورودی مدل‌های پیشنهاد شده در این معماری داده‌های چندحالتی<sup>۳۱</sup> هستند که به صورت همزمان و یکجا به مدل داده می‌شوند برای مثال تصویر به همراه یک جمله توصیف‌کننده آن و یا یک فیلم به همراه زیرنویسش به این شبکه‌ها برای آموزش داده می‌شوند. به علاوه این مدل‌ها با ترکیبی از اهداف مختلف مانند visual-based Masked Language Model، text-based Masked Language Model، - masked visual-

<sup>۲۵</sup> vision-and-language pretraining models

<sup>۲۶</sup> single-stream

<sup>۲۷</sup> two-stream

<sup>۲۸</sup> text encoding

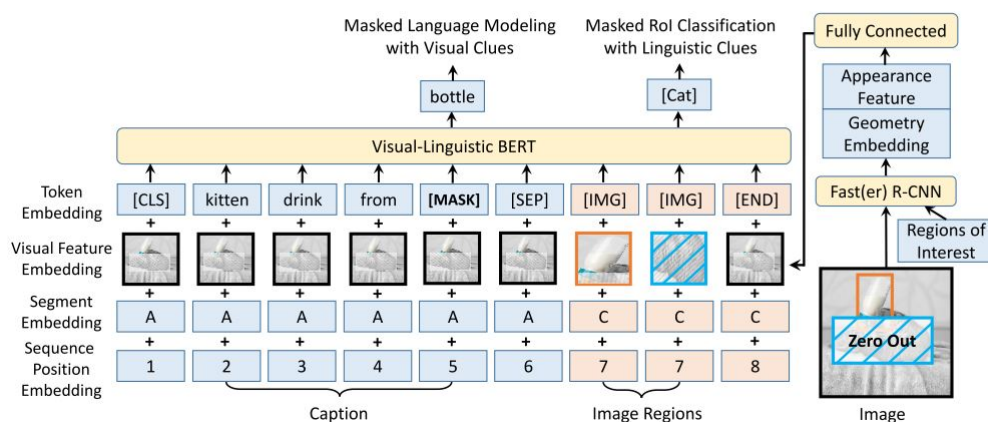
<sup>۲۹</sup> image encoding

<sup>۳۰</sup> encoder

<sup>۳۱</sup> multimodal



فصل ۳. مروری بر کارهای مرتبط ۳-۴. رویکرد مدل‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر



شکل ۳-۱۲: معماری شبکه از قبل آموزش دیده VL-BERT [۶۱]

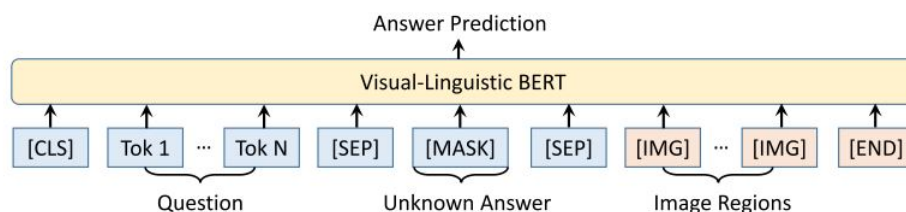
feature modeling و visual-linguistic matching بهینه می‌شود. سپس از بازنمایی‌های آموخته شده توسط این مدل‌ها در مسائل پایین دستی understanding و یا generation استفاده می‌شود. به عنوان مثال، مدل VideoBERT [۶۲] برای مسائل generation مانند تولید توصیف فیلم طراحی شده است. در حالی که چندین مدل دیگر مانند B2T2 [۲]، Unicoder-VL [۳۳]، VL-BERT [۶۱] و UNITER [۱۰] وجود دارد که همگی برای مسائل understanding طراحی شده‌اند. مدل‌های دیگری مانند VLP [۷۹] و OSCAR [۳۴] مدل‌های یکپارچه‌ای هستند که هم در مسائل پایین دستی understanding و هم در مسائل generative کاربرد دارد. از بین این مدل‌ها، تنها از مدل‌های VL-BERT، UNITER، VLP و OSCAR می‌توان برای مسئله پرسش و پاسخ تصویری استفاده کرد. بنابراین در ادامه این بخش جزئیات هر کدام از این مدل‌ها را توضیح خواهیم داد.

### ۳-۴-۱-۱ شبکه VL-BERT

شکل ۳-۱۲ معماری VL-BERT را نشان می‌دهد. مشابه BERT، از کدگذارهای multi-layer bidirectional transformer استفاده شده است. اما برخلاف BERT که ورودی آن تنها کلمات جمله هستند، این شبکه به همراه کلمات یک جمله، مناطق مورد علاقه<sup>۳۲</sup> استخراج شده از تصویر و یا به اختصار ROI را نیز به عنوان ورودی می‌گیرد. برای استخراج ROI از تصویر از شبکه Faster RCNN [۵۳] استفاده شده است. هر ورودی این شبکه با توکن [CLS] آغاز می‌شود. سپس با کلمات جمله و ROI های تصویر ادامه می‌یابد و با توکن

<sup>۳۲</sup>regions-of-interest

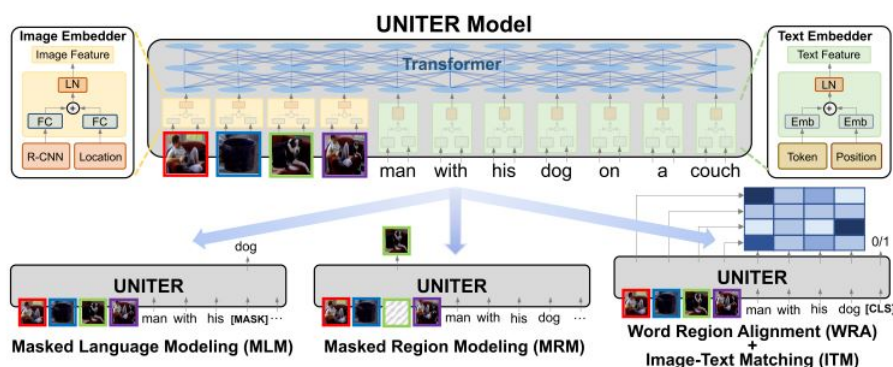
فصل ۳. مروری بر کارهای مرتبط ۳-۴. رویکرد مدل‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر



شکل ۳-۱۳: نحوه ورودی و خروجی شبکه VL-BERT برای آموزش در مسئله پرسش و پاسخ تصویری [۶۱]

[END] خاتمه می‌یابد. از توکن [SEP] نیز برای جدا کردن جملات و یا جملات و تصویر از هم استفاده می‌شود. برای هر ورودی، تعبیه ویژگی<sup>۳۳</sup> آن جمع چهار نوع تعبیه است که در شکل ۳-۱۲ مشخص شده است. در میان آن‌ها، تعبیه مربوط به ویژگی‌های تصویری<sup>۳۴</sup> به تازگی به شبکه اضافه شده است در حالی که سه تعبیه دیگر از قبل در مدل BERT وجود داشته است. برای آموزش VL-BERT از دادگان Conceptual Captions به عنوان دادگان زبانی- تصویری استفاده شده است. علاوه بر این از دو دادگان فقط زبانی به نام‌های English Wikipedia و BooksCorpus به منظور بهبود تعمیم‌دهی شبکه استفاده شده است. برای بهینه‌سازی شبکه VL-BERT از دو تابع هدف استفاده شده است: text-based Masked Language Model و visual-based Masked Language Model. در text-based Masked Language Model با احتمال ۱۵ درصد یکی از کلمات ورودی با توکن [MASK] جایگزین می‌شود. بنابراین شبکه باید سعی کند که این کلمه ماسک شده را با توجه به کلمات دیگر و ویژگی‌های تصویری در خروجی پیش‌بینی نماید. در visual-based Masked Language Model با احتمال ۱۵ درصد یکی از ROI ها ماسک می‌شود و شبکه باید سعی کند در خروجی برچسب گروه مربوط به آن ROI را با توجه به کلمات و سایر ROI ها پیش‌بینی کند. دقت شود که همانطور که در قسمت سمت راست تصویر ۳-۱۲ مشخص است، ملاک برچسب گروه‌بندی درست برای ROI ها، خروجی شبکه Faster RCNN است. برای استفاده از شبکه از قبل آموزش دیده VL-BERT برای مسئله پرسش و پاسخ تصویری، مطابق شکل ۳-۱۳ سه تایی کلمات سوال، پاسخ و ROI های استخراج شده از تصویر توسط Faster RCNN در ورودی داده می‌شود که به جای پاسخ، [MASK] قرار گرفته که شبکه تلاش می‌کند؛ پاسخ را در خروجی پیش‌بینی کند.

فصل ۳. مروری بر کارهای مرتبط ۳-۴. رویکرد مدل‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر



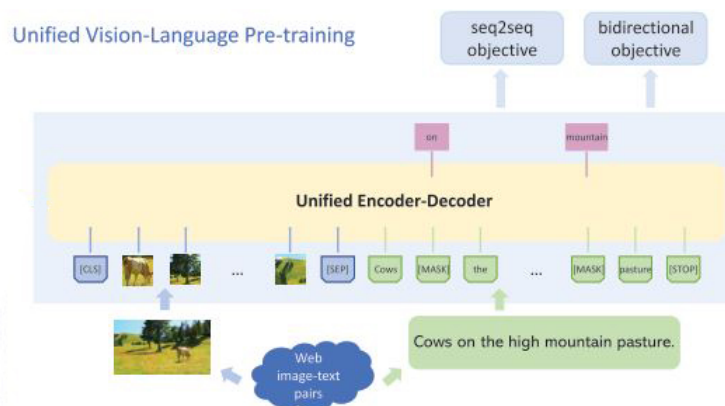
شکل ۳-۱۴: معماری شبکه از قبل آموزش دیده UNITER [۱۰]

### ۳-۴-۱-۲ شبکه UNITER

معماری مدل UNITER در شکل ۳-۱۴ نشان داده شده است. ورودی این مدل مانند VL-BERT، کلمات یک جمله به همراه ROI های تصویر است. یکی از تفاوت‌های مدل UNITER با مدل VL-BERT این است که از ۴ دادگان زبانی-تصویری برای آموزش استفاده کرده است: (۱) COCO، (۲) Visual Genome، (۳) Conceptual Captions و (۴) SBU Captions. تفاوت دیگر این مدل با مدل VL-BERT در توابع هدف است که علاوه بر text-based Masked Language Model و visual-based Masked Language Model از دو تابع هدف دیگر به نام‌های Image-Text Matching و Word-Region Alignment نیز استفاده می‌کند. در Image-Text Matching هدف این است که مدل بتواند پیش‌بینی کند که آیا جمله و تصویر داده شده در ورودی با هم مطابقت دارند یا خیر. بدین منظور، یک جمله و ROI های تصویر به UNITER داده می‌شود و در خروجی بازنمایی مربوط به توکن [CLS] از یک تابع سیگموئید عبور داده می‌شود که یک مقدار بین صفر و یک را برمی‌گرداند که مقدار یک نشان می‌دهد که جمله و تصویر ورودی کاملاً با هم مطابقت دارد و مقدار صفر به این معناست که جمله و تصویر ورودی با هم مطابقت ندارد. در UNITER علاوه بر در نظر گرفتن تطابق جمله و تصویر، از تطابق بین کلمات موجود در جمله و ROI های تصویر نیز برای آموزش استفاده می‌شود که این موضوع در قالب تابع هدف Word-Region Alignment در مدل مطرح شده است. زمان آموزش مدل UNITER به ازای هر دسته از داده‌های ورودی، یکی از ۴ تابع هدف نامبرده شده به صورت تصادفی انتخاب می‌شود و براساس آن تابع هدف، عملیات کاهش گرادیان برای شبکه انجام می‌شود. برای استفاده از شبکه

<sup>۳۳</sup>feature embedding  
<sup>۳۴</sup>visual feature embedding

فصل ۳. مروری بر کارهای مرتبط ۳-۴. رویکرد مدل‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر



شکل ۳-۱۵: معماری شبکه از قبل آموزش دیده VLP [۷۹]

از قبل آموزش دیده UNITER برای مسئله پرسش و پاسخ تصویری، بازنمایی حاصل از توکن [CLS] به یک شبکه MLP داده می‌شود و پاسخ را برای سوال و تصویر ورودی پیش‌بینی می‌کند. در واقع در این حالت، مسئله پرسش و پاسخ تصویری به عنوان یک مسئله طبقه‌بندی در نظر گرفته می‌شود.

### ۳-۴-۱-۳ شبکه VLP

شبکه از قبل آموزش دیده VLP نیز مانند دو شبکه‌ی قبلی از کلمات یک جمله و ROI های استخراج شده از تصویر به عنوان ورودی استفاده می‌کند. تفاوت اصلی این شبکه با دو شبکه VL-BERT و UNITER در این است که یک شبکه‌ی یکپارچه رمزگذار-رمزگشا است که نه تنها در مسائل understanding بلکه در مسائل generative به دلیل وجود رمزگشا قابل استفاده است. مدل VLP بر روی دادگان Conceptual Captions آموزش داده شده است. دو تابع هدف در شبکه VLP استفاده شده است: (۱) bidirectional و seq2seq. در تابع هدف bidirectional یکی از کلمات موجود در جمله با توکن [MASK] جایگزین می‌شود و برای پیش‌بینی این کلمه ماسک شده در خروجی از تمامی کلمات و ROI های اطراف آن استفاده می‌شود. اما در تابع هدف seq2seq برای پیش‌بینی کلمه ماسک شده در خروجی، تنها از کلمات سمت چپ کلمه ماسک شده و ROI های اطراف آن استفاده می‌شود. به عبارتی دیگر، برای پیش‌بینی کلمه ماسک شده نمی‌توان از کلماتی که بعد از آن و در آینده در جمله آمده است؛ استفاده کرد. معماری شبکه VLP در شکل ۳-۱۵ نشان داده شده است.

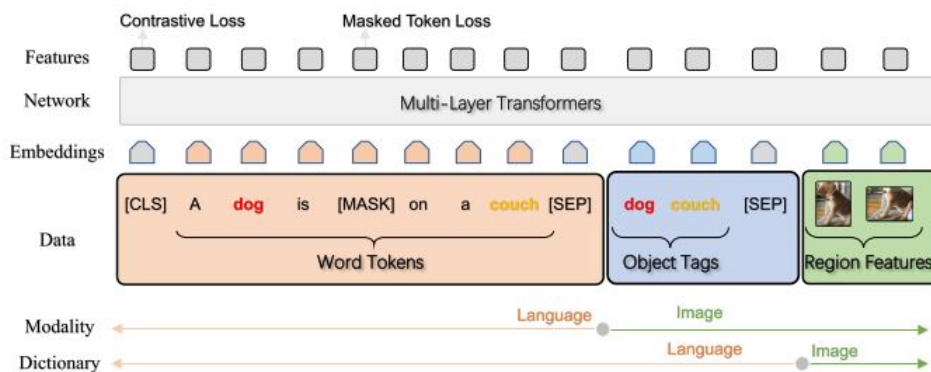
فصل ۳. مروری بر کارهای مرتبط ۳-۴. رویکرد مدل‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر

### ۴-۱-۴-۳ شبکه OSCAR

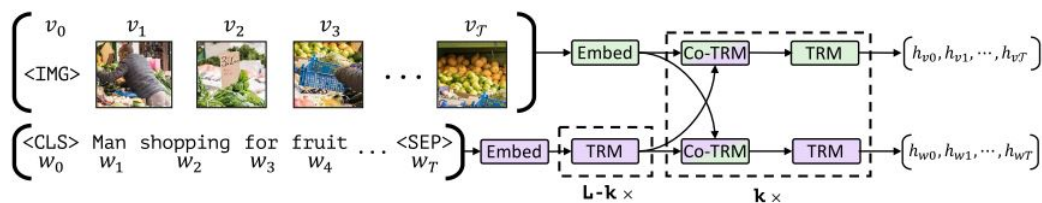
ورودی سه مدل قبلی یعنی VL-BERT ، UNITER و VLP یک جمله به همراه ROI های استخراج شده از تصویر بود. در مدل OSCAR علاوه بر این دو ورودی از ورودی دیگری به نام برچسب اشیا<sup>۳۵</sup> استفاده می‌شود که اشیایی که هم در تصویر وجود دارد و هم در جمله به آن اشاره شده است را نشان می‌دهد. در [۲۴] ادعا شده است که استفاده برچسب اشیا منجر به تولید بازنمایی بهتری از متن و تصویر می‌شود و در واقع از این برچسب‌ها به عنوان لنگر برای تطابق دادن فضای تصویر و متن استفاده می‌شود. در مدل OSCAR برای بدست آوردن ROI های تصویر و برچسب اشیا از شبکه Faster RCNN استفاده شده است. در مدل OSCAR به دو طریق می‌توان به ورودی‌ها نگاه کرد که در نتیجه دو تابع هدف برای آموزش این شبکه تعریف می‌شود. در روش اول، کلمات جمله و برچسب اشیا با هم در نظر گرفته می‌شود (دید Dictionary) و به احتمال ۱۵ درصد یکی از کلمات جمله و یا یکی از برچسب‌های اشیا با توکن [MASK] جایگزین می‌شود و مدل باید سعی کند این کلمه ماسک شده را در خروجی پیش‌بینی کند (Masked Token Loss). در روش دوم، ROI های تصویر و برچسب اشیا با هم در نظر گرفته می‌شود (دید Modality) و با احتمال ۵۰ درصد برچسب‌های اشیا با برچسب‌های دیگری تغییر می‌کند و مدل باید پیش‌بینی کند که آیا کلمات موجود در جمله با قسمت برچسب اشیا و ROI های تصویر مطابقت دارد یا نه. که بدین منظور خروجی شبکه برای توکن [CLS] به یک شبکه کاملاً متصل داده می‌شود و یک طبقه‌بندی باینری انجام می‌شود که یک به معنای تطابق کلمات جمله با ROI های تصویر و برچسب اشیاست و صفر نشان‌دهنده عدم تطابق است (Contrastive Loss). برای آموزش مدل OSCAR از مجموعه داده‌های COCO ، Conceptual Captions ، SBU captions ، flicker30 و GQA استفاده شده است. برای استفاده از شبکه از قبل آموزش دیده OSCAR برای مسئله پرسش و پاسخ تصویری، سوال به همراه برچسب اشیا و ROI های تصویر به ورودی شبکه داده می‌شود و خروجی توکن [CLS] به یک طبقه‌بند داده می‌شود تا پاسخ سوال و تصویر داده شده در تصویر بدست آید. در واقع در این روش، مسئله پرسش و پاسخ تصویری به صورت یک مسئله طبقه‌بندی در نظر گرفته می‌شود. معماری شبکه OSCAR در شکل ۳-۱۶ نمایش داده شده است.

<sup>۳۵</sup>object tag

فصل ۳. مروری بر کارهای مرتبط ۳-۴. رویکرد مدل‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر



شکل ۳-۱۶: معماری شبکه از قبل آموزش دیده OSCAR [۳۴]



شکل ۳-۱۷: معماری شبکه از قبل آموزش دیده ViLBERT [۳۸]

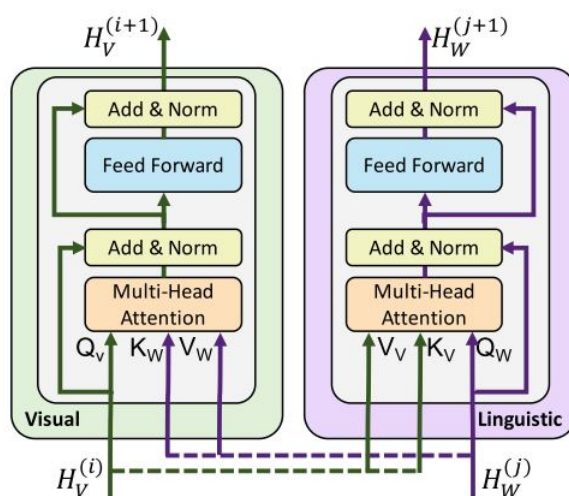
### ۳-۴-۲ معماری دو جریان

در مقابل معماری تک جریان، معماری دو جریان برای یادگیری هر کدام از بازنمایی‌های تصویر و متن از یک رمزگذار مستقل استفاده می‌کند. سپس از یک رمزگذار دیگر برای بدست آوردن بازنمایی مشترک متن و تصویر استفاده می‌کند. مشابه معماری تک جریان، معماری دو جریان نیز مدل‌های خود را با visual-based Masked Language Model، text-based Masked Language Model و visual-linguistic matching بهینه می‌کنند. ViLBERT [۳۸] و LXMERT [۶۵] نمونه‌هایی از معماری دو جریان هستند که از این دو مدل می‌توان برای مسئله پرسش و پاسخ تصویری استفاده کرد. پس در ادامه این بخش، جزئیات این دو شبکه را بررسی خواهیم کرد.

### ۳-۴-۱-۲ شبکه ViLBERT

شکل ۳-۱۷ معماری شبکه ViLBERT را نمایش می‌دهد. مدل ViLBERT شامل دو مدل موازی به سبک BERT است که به صورت جداگانه بر روی کلمات متن و ROI های تصویر اعمال می‌شود و از بلوک‌های

فصل ۳. مروری بر کارهای مرتبط ۳-۴. رویکرد مدل‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر



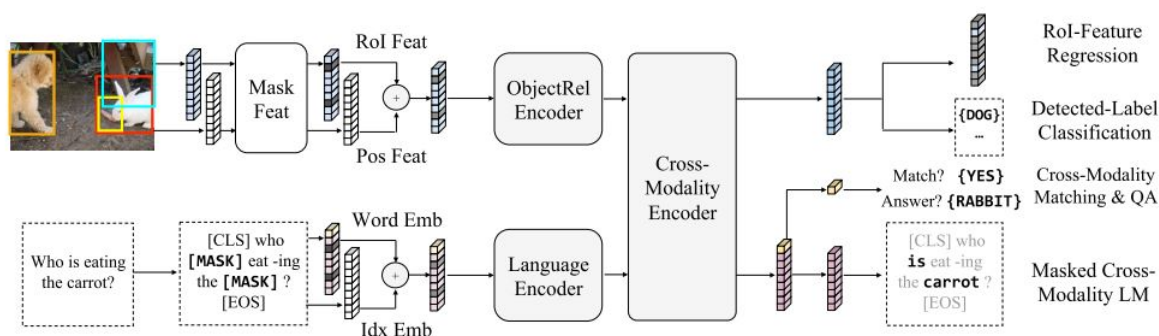
شکل ۳-۱۸: ساختار لایه co-attentional transformer [۳۸]

ترنسفرمر در هر جریان استفاده شده است (در شکل ۳-۱۷ با TRM مشخص شده است). سپس برای بدست آوردن بازنمایی مشترک بین متن و تصویر از لایه‌های co-attentional transformer استفاده شده است (در شکل ۳-۱۷ با Co-TRM مشخص شده است). اساس لایه co-attentional transformer بر پایه‌ی ترنسفرمر است در واقع برای هر کدام از بخش‌های تصویری و متنی داده ورودی، یک ترنسفرمر در لایه co-attentional transformer در نظر گرفته شده است که پس از عبور متن و داده از جریان‌های مستقل خود و بدست آمدن query، key و value برای هر کدام، key و value متن به ترنسفرمر تصویر در co-attentional transformer داده می‌شود و به صورت متقابل key و value تصویر به ترنسفرمر متن داده می‌شود.

شکل ۳-۱۸ ساختار لایه co-attentional transformer را نشان می‌دهد. برای آموزش مدل ViLBERT از توابع هدف visual-based Masked Language Model، text-based Masked Language Model و visual-linguistic matching استفاده شده است. شبکه ViLBERT بر روی دادگان Conceptual Captions آموزش داده شده است. برای استفاده از شبکه از قبل آموزش دیده ViLBERT برای مسئله پرسش و پاسخ تصویری، ابتدا خروجی بازنمایی توکن [CLS] و بازنمایی تصویر ضرب متناظر می‌شوند. سپس با عبور از یک شبکه MLP دولایه پاسخ مربوط به سوال و تصویر حاصل می‌شود.



فصل ۳. مروری بر کارهای مرتبط ۳-۴. رویکرد مدل‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر



شکل ۳-۱۹: معماری شبکه از قبل آموزش دیده LXMERT [۶۵]

### ۳-۴-۲ شبکه LXMERT

شکل ۳-۱۹ معماری مدل LXMERT را نشان می‌دهد. ورودی این شبکه کلمات جمله ورودی و ROI های استخراج شده از تصویر است. همان‌طور که قبلاً اشاره شد؛ مدل LXMERT یک مدل دو جریان است به همین دلیل برای پردازش متن و تصویر از دو رمزگذار مجزا و مستقل استفاده شده است (در شکل ۳-۱۹ به ترتیب با عنوان‌های ObjectRel Encoder و Language Encoder برای تصویر و متن مشخص شده است. ( و سپس برای بدست آوردن بازنمایی مشترک از رمزگذار Cross-modality استفاده شده است. توابع هدف استفاده شده در مدل LXMERT مشابه شبکه ViLBERT است اما در LXMERT از تابع هدف دیگری به نام image question answering برای آموزش شبکه استفاده شده است. زیرا حدود ۱/۳ داده‌ای که برای آموزش این شبکه استفاده شده است؛ یک سوال در مورد تصویر ورودی است. بنابراین با تعریف تابع هدف image question answering مدل سعی می‌کند تا پاسخ این سوال را در خروجی پیش‌بینی کند. برای آموزش شبکه LXMERT از مجموعه داده‌های MS COCO ، Visual Genome ، VQA v2.0 ، GQA balanced version و VG-QA استفاده شده است.

در جدول ۳-۶ مقایسه چند نمونه از مدل‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر که مسئله پرسش و پاسخ تصویری را پشتیبانی می‌کنند؛ آورده شده است. ورودی تمام این مدل‌ها، کلمات جمله و ROI های تصویر است به جز مدل OSCAR که علاوه بر این دو، برچسب اشیا را نیز به عنوان ورودی دریافت می‌کند. شباهت دیگر این مدل‌ها در استفاده از دادگان Conceptual Captions برای آموزش است البته به جز مدل LXMERT که از این دادگان استفاده نکرده است. نکته‌ی حائز اهمیت دیگر در این جدول استفاده تقریباً



فصل ۳. مروری بر کارهای مرتبط ۳-۴. رویکرد مدل‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر

جدول ۳-۶: مقایسه بین شبکه‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر

روش	معماری	ورودی	مجموعه داده‌گان استفاده شده برای آموزش	توابع هدف
<b>VL-BERT</b> [۶۱]	تک جریان	کلمات جمله ROI‌های تصویر	Conceptual Captions + BooksCorpus + English Wikipedia	text-based MLM + visual-based MLM
<b>UNITER</b> [۱۰]	تک جریان	کلمات جمله ROI‌های تصویر	COCO + Visual Genome + Conceptual Captions + SBU Captions	text-based MLM + visual-based MLM + Image-Text Matching + Word- Region Alignment
<b>VLP</b> [۷۹]	تک جریان	کلمات جمله ROI‌های تصویر	Conceptual Captions	bidirectional + seq2seq
<b>OSCAR</b> [۳۴]	تک جریان	کلمات جمله ROI‌های تصویر + برچسب اشیا	COCO + Conceptual Captions + SBU captions + flicker30 + GQA	Masked Token Loss + Contrastive Loss
<b>ViL-BERT</b> [۳۸]	دو جریان	کلمات جمله ROI‌های تصویر	Conceptual Captions	text-based MLM + visual-based MLM + Image- Text Matching
<b>LXMERT</b> [۶۵]	دو جریان	کلمات جمله ROI‌های تصویر	MS COCO + Visual Genome + VQA v2.0 + GQA balanced version + VG-QA	text-based MLM + visual-based MLM + Image- Text Matching + Image Question Answering

جدول ۳-۷: دقت شبکه‌های از قبل آموزش دیده بر روی مجموعه داده VQA v2.0 (test-std)

روش	سوالات بله/خیر	سوالات شمارشی	سایر سوالات	دقت کل
VLP[۷۹]	۸۷/۴	۵۲/۱	۶۰/۵	۷۰/۷
ViL-BERT[۳۸]	—	—	—	۷۰/۹۲
VL-BERT[۶۱]	—	—	—	۷۲/۲۲
LXMERT[۶۵]	۸۸/۲	۵۴/۲	۶۳/۱	۷۲/۵
OSCAR[۳۴]	—	—	—	۷۳/۸۲
UNITER[۱۰]	—	—	—	۷۴/۰۲

تمامی مدل‌ها از دو تابع هدف text-based Masked Language Model و visual-based Masked Language Model است.

در جدول ۳-۷ نتایج مدل‌های ViL-BERT، OSCAR، VLP، UNITER، VL-BERT و LXMERT بر روی مجموعه داده VQA v2.0 نشان داده شده است. بهترین نتیجه بدست آمده برای مدل UNITER است. یکی از نکات قابل ملاحظه در این جدول این است که مدل‌های تک جریان نتایج بهتری نسبت به مدل‌های دو جریان بدست آوردند در حالی که تعداد پارامترهای مدل‌های تک جریان نسبت به مدل‌های دو جریان کمتر است.

### ۳-۵ معیارهای ارزیابی مسئله پرسش و پاسخ تصویری

در این بخش می‌خواهیم به طور مختصر معیارهای ارزیابی شناخته شده در مسئله پرسش و پاسخ تصویری را بررسی کنیم. همان‌طور که قبلاً ذکر شد؛ معمولاً دو نوع سوال در مجموعه داده‌های پرسش و پاسخ تصویری در نظر گرفته می‌شود: سوالات open-ended و سوالات چندگزینه‌ای. در سوالات چندگزینه‌ای، برای هر سوال دقیقاً یک پاسخ صحیح وجود دارد. بنابراین ارزیابی آن ساده است زیرا می‌توان به راحتی از معیار دقت استفاده کرد. اما در سوالات open-ended این امکان وجود دارد که چندین پاسخ صحیح برای هر سوال وجود داشته باشد. بنابراین ارزیابی در این حالت ساده نخواهد بود. برای حل این موضوع، اکثر مجموعه داده‌های پرسش و پاسخ تصویری پاسخ‌ها را محدود به چند کلمه (۱ تا ۳ کلمه) می‌کنند و یا پاسخ‌ها را از یک مجموعه بسته انتخاب می‌کنند. در ادامه به بررسی مهم‌ترین معیارهای این حوزه می‌پردازیم. اما ارزیابی مسئله پرسش و پاسخ تصویری همچنان یک مسئله حل نشده است. هر کدام از روش‌ها و معیارهای ارزیابی موجود، مزیت‌ها و

معایب خاص خود را دارند. بنابراین برای انتخاب معیار ارزیابی باید به مواردی همچون ساختار دادگان و نحوه ساخت آن، میزان بایاس موجود در دادگان توجه نمود.

### ۳-۵-۱ معیار دقت

اگر چه در سوالات چندگزینه‌ای برای سنجش یک مدل معیار دقت کافی است اما در سوالات open-ended معیار دقت سخت‌گیرانه است زیرا فقط در حالتی که پاسخ مدل کاملاً مطابق با پاسخ در نظر گرفته شده باشد، پذیرفته می‌شود. برای مثال اگر صورت سوال «چه حیواناتی در تصویر است؟» باشد و پاسخ مدل به جای «سگ‌ها» پاسخ «سگ» باشد؛ غلط تلقی می‌شود. بنابراین به دلیل این محدودیت‌هایی که معیار دقت دارد؛ معیارهای دیگری برای ارزیابی این نوع سوالات پیشنهاد شده‌است.

$$Accuracy = \frac{\text{Number of questions answered correctly}}{\text{Total questions}} \quad (2-3)$$

### ۳-۵-۲ معیار شباهت Wu-Palmer [۷۳]

این معیار ارزیابی توسط مالینوفسکی [۴۲] برای پرسش و پاسخ تصویری ارائه شد. این معیار از تئوری مجموعه‌های فازی الهام گرفته شده است و نسبت به معیار دقت سخت‌گیری کمتری دارد. معیار شباهت Wu-Palmer سعی می‌کند که تفاوت بین پاسخ پیش‌بینی شده با پاسخ صحیح را از لحاظ معنایی اندازه‌گیری کند. یکی از معایب این معیار این است که به پاسخ‌هایی که از لحاظ لغوی شبیه هم هستند ولی از لحاظ معنایی متفاوت هستند، امتیاز بالایی می‌دهد. زمانی که پاسخ‌های ما به صورت عبارت یا جمله باشد؛ این معیار عملکرد خوبی ندارد.

### ۳-۵-۳ معیار اجماع

از این معیار زمانی استفاده می‌شود که هر سوال توسط کاربرهای انسانی متفاوتی پاسخ داده شود. در واقع برای هر سوال چندین پاسخ مستقل وجود داشته باشد. این معیار دو نوع دارد: میانگین اجماع و کمترین اجماع. در میانگین اجماع امتیاز نهایی برابر با میانگین وزندار پاسخ‌های وارد شده توسط کاربرهای متفاوت است و در کمترین اجماع پاسخ پیش‌بینی شده حداقل باید با یکی از پاسخ‌ها مطابقت داشته باشد. در مسئله‌ی

پرسش و پاسخ تصویری معمولاً از حالت کمترین اجماع استفاده می‌شود و آستانه را هم برابر ۳ قرار می‌دهند به این معنی که اگر پاسخ پیش‌بینی شده با ۳ یا بیشتر از ۳ پاسخ برابر باشد امتیاز کامل می‌گیرد و در غیر این صورت هیچ امتیازی کسب نخواهد کرد. از معایب این روش می‌توان به هزینه زیاد جمع‌آوری پاسخ برای سوالات اشاره کرد. آنتول و همکارانش از این معیار ارزیابی در [۴] استفاده کرده‌اند.

$$Accuracy_{VQA} = \min\left(\frac{n}{3}, 1\right) \quad (3-3)$$

### ۳-۵-۴ معیار MPT

یکی از مشکلات مجموعه داده‌های پرسش و پاسخ تصویری توزیع غیریکنواخت انواع سوال‌هاست. در این مواقع، نمی‌توان از معیار دقت استفاده کرد. بنابراین در [۲۷] معیار جدیدی به نام MPT<sup>۳۶</sup> ارائه شده است که توزیع نامتوازن سوال‌ها را جبران می‌کند. معیار MPT میانگین دقت برای هر نوع سوال را محاسبه می‌کند. از نسخه‌ی نرمالایز شده‌ی این معیار نیز برای رفع مشکل بایاس در توزیع پاسخ‌ها استفاده می‌شود.

### ۳-۵-۵ معیار BLEU

معیار BLEU [۴۷]<sup>۳۷</sup> یکی از معیارهای ارزیابی خودکار ترجمه ماشینی است. در [۲۰] پیشنهاد داده شد که از این معیار نیز برای ارزیابی پرسش و پاسخ تصویری می‌توان استفاده کرد. معیار BLEU کنار هم قرار گرفتن n-gram های پاسخ پیش‌بینی شده و پاسخ صحیح را اندازه‌گیری می‌کند. معمولاً BLEU زمانی که جمله‌ها کوتاه باشند، با شکست مواجه می‌شود.

### ۳-۵-۶ معیار METEOR

معیار METEOR [۱۴]<sup>۳۸</sup> نیز همانند BLEU یکی از معیارهای ارزیابی خودکار ترجمه ماشینی است. به پیشنهاد [۲۰] از این معیار هم می‌توان برای پرسش و پاسخ تصویری نیز استفاده نمود. معیار METEOR سعی می‌کند که هم‌ترازی بین کلمات موجود در پاسخ پیش‌بینی شده و پاسخ صحیح را پیدا کند.

<sup>۳۶</sup>Mean Per Type

<sup>۳۷</sup>BiLingual Evaluation Understudy

<sup>۳۸</sup>Metric for Evaluation of Translation with Explicit ORdering

## ۳-۶ جمع‌بندی

در این فصل، پس از مقایسه مجموعه‌دادگان مختلف در پرسش و پاسخ تصویری، به سراغ رویکردهای حل این مسئله از دو منظر یادگیری عمیق و شبکه‌های از قبل آموزش دیده بر روی زبان طبیعی و تصویر رفتیم. عموماً روش‌هایی که در رویکرد یادگیری عمیق پیشنهاد شده‌اند؛ دارای سه فاز هستند. در فاز اول از تصویر و سوال ویژگی استخراج می‌شود و در فاز دوم از روش‌های ساده مانند ضرب ویژگی‌ها تا روش‌های پیچیده‌تر مانند مکانیزم توجه استفاده می‌شود تا بازنمایی مشترک بین تصویر و سوال بدست آید. در فاز آخر از این بازنمایی مشترک برای بدست آوردن پاسخ در خروجی استفاده می‌شود. در رویکرد شبکه‌های از قبل آموزش دیده، براساس نحوه کدگذاری متن و تصویر که به صورت همزمان یا موازی انجام شود؛ شبکه‌ها را به دو معماری تک جریان و دو جریان تقسیم کردیم. برای هر کدام از معماری‌ها چند نمونه را معرفی و بررسی کردیم. در انتهای این فصل هم به شرح معیارهای ارزیابی مسئله پرسش و پاسخ تصویری پرداختیم.

## فصل ۴

# نتیجه‌گیری و کارهای آینده

### ۴-۱ نتیجه‌گیری

علی‌رغم این که از معرفی مسئله پرسش و پاسخ تصویری تنها چندین سال می‌گذرد، رشد آن در این چند سال قابل توجه بوده است. مجموعه‌دادگان بسیاری با اهداف مختلف در طی این سال‌ها معرفی شد. برای حل مسئله پرسش و پاسخ تصویری، رویکردهای یادگیری عمیق همچنان در مرکز توجه هستند. ما برجسته‌ترین مدل‌های یادگیری عمیق برای مسئله پرسش و پاسخ تصویری را بررسی کردیم. با معرفی شبکه‌های از قبل آموزش‌دیده، بهبود چشمگیری در مسائل یادگیری عمیق رخ داد به طوری که بیشتر مسائل مختلف در یادگیری عمیق، بهترین نتیجه خود را با استفاده از شبکه‌های از قبل آموزش‌دیده بدست آورده‌اند. مسئله پرسش و پاسخ تصویری نیز از این قاعده مستثنی نیست و در حال حاضر شبکه‌های از قبل آموزش‌دیده بر روی زبان طبیعی و تصویر بهترین عملکرد را برای مجموعه‌دادگان پرسش و پاسخ تصویری رقم زده‌اند. چندین نمونه از این مدل‌ها را با جزئیات بحث کردیم. در آخر معیارهایی را معرفی کردیم که بتوان با آن‌ها مدل‌های پرسش و پاسخ تصویری را ارزیابی کرد. البته که ارزیابی مسئله پرسش و پاسخ تصویری همچنان یک مسئله حل نشده است و نیاز به تحقیقات بیشتری دارد. پیشرفت‌های زیادی که همچنان برای مجموعه‌دادگان مختلف در این حوزه اتفاق می‌افتد، به این معناست که هنوز فضای زیادی برای نوآوری در آینده وجود دارد.

## ۴-۲ مسائل باز و کارهای قابل انجام

با وجود تمام پیشرفت‌هایی که در سال‌های اخیر در مسئله پرسش و پاسخ تصویری اتفاق افتاده است، مدل‌های پیشنهاد شده در این حوزه با نواقصی مواجه هستند. اولین مشکل روش‌های فعلی پاسخ به سوالاتی است که نیاز به استدلال طولانی دارند. از طرفی منبع بهبودهای نسبی مدل‌های موجود واضح نیست و مشخص نیست که مدل تا چه اندازه مفاهیم مشترک بین زبان و تصویر را درک می‌کند و چگونه از پیوند این دو برای پیش‌بینی پاسخ استفاده می‌کند. پس اگر بتوانیم بفهمیم که روند درک مدل‌هایی فعلی از زبان و تصویر چگونه است، می‌توانیم مدلی را پیشنهاد دهیم که بتواند به سوالاتی که نیاز به استدلال طولانی دارند، پاسخ دهد.

اکثر روش‌های پیشنهادشده، مسئله پرسش و پاسخ تصویری را یک مسئله طبقه‌بندی در نظر می‌گیرند و تعداد کمی از کارهای انجام شده به دنبال تولید پاسخ بوده‌اند. یکی از دلایلی که باعث کم توجهی به تولید پاسخ شده است، زمان‌بر بودن فرآیند آن است. یکی از راه‌حل‌های این مشکل می‌تواند استفاده از ترنسفرمرها با چندین لایه رمزگذار و رمزگشا بر روی هم باشد. از معماری ترنسفرمر برای تولید پاسخ در پرسش و پاسخ تصویری به صورت محدود استفاده شده است. از طرفی، موفقیت ترنسفرمرها در مسائل پردازش زبان طبیعی، ما را ترغیب می‌کند که از قدرت آن‌ها در مسئله پرسش و پاسخ تصویری برای تولید پاسخ در آینده استفاده کنیم.

یکی دیگر از محدودیت‌های مسئله پرسش و پاسخ تصویری، فقدان مجموعه‌دادگان متناسب با واقعیت است. در حال حاضر نمی‌توان از دادگان موجود در مسئله پرسش و پاسخ تصویری برای کاربردهای عملی مانند کمک به افراد نابینا و کم‌بینا استفاده کرد. از طرف دیگر اکثر مجموعه‌دادگان با مشکل بایاس مواجه هستند. بنابراین جمع‌آوری و تهیه مجموعه‌دادگانی که منطبق با کاربرد عملی در جامعه و بدون بایاس باشند، اهمیت پیدا می‌کند.

در فصل قبل دیدیم که در حال حاضر، بهترین عملکرد برای مجموعه‌دادگان پرسش و پاسخ تصویری توسط شبکه‌های از قبل آموزش‌دیده بر روی زبان طبیعی و تصویر بدست آمده است. اساس و پایه‌ی این شبکه‌ها، ترنسفرمر است. یکی از بزرگترین مشکلات ترنسفرمرها این است که محاسبه توجه از مرتبه زمانی و حافظه‌ای ۲ است. اخیراً روش‌های زیادی مانند Reformer [۳۰] و Performer [۱۲] پیشنهاد شده است که مرتبه زمانی و حافظه‌ای ترنسفرمرها را کاهش می‌دهند. بنابراین یکی از مسیرهای تحقیقاتی پیش رو، استفاده از این ترنسفرمرهای بهبودیافته در معماری شبکه‌های از قبل آموزش‌دیده بر روی زبان طبیعی و تصویر می‌تواند

باشد.

با توجه به دانشی که ما بدست آوردیم، تاکنون هیچ‌گونه تحقیقی در مورد پرسش و پاسخ تصویری در زبان فارسی انجام نشده است. از این رو دادگان مناسبی نیز برای این کار وجود ندارد. پس تهیه و جمع‌آوری دادگان فارسی برای مسئله پرسش و پاسخ تصویری و آموزش یک مدل کارآمد براساس آن، یک کار ارزشمند خواهد بود و مسیر جدیدی را برای سایر محققین باز خواهد کرد.



- [1] ACHARYA, M., KAFLE, K., AND KANAN, C. Tallyqa: Answering complex counting questions. in *Proceedings of the AAAI Conference on Artificial Intelligence* (2019), volume 33, pp. 8076–8084.
- [2] ALBERTI, C., LING, J., COLLINS, M., AND REITTER, D. Fusion of detected objects in text for visual question answering. in *EMNLP/IJCNLP* (2019).
- [3] ANDREAS, J., ROHRBACH, M., DARRELL, T., AND KLEIN, D. Neural module networks. in *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016), pp. 39–48.
- [4] ANTOL, S., AGRAWAL, A., LU, J., MITCHELL, M., BATRA, D., LAWRENCE ZITNICK, C., AND PARIKH, D. Vqa: Visual question answering. in *Proceedings of the IEEE international conference on computer vision* (2015), pp. 2425–2433.
- [5] BAI, Y., FU, J., ZHAO, T., AND MEI, T. Deep attention neural tensor network for visual question answering. in *Proceedings of the European Conference on Computer Vision (ECCV)* (2018), pp. 20–35.
- [6] BEN-YOUNES, H., CADENE, R., CORD, M., AND THOME, N. Mutan: Multimodal tucker fusion for visual question answering. in *Proceedings of the IEEE international conference on computer vision* (2017), pp. 2612–2620.
- [7] BROWN, T. B., MANN, B., RYDER, N., SUBBIAH, M., KAPLAN, J., DHARIWAL, P., NEELAKANTAN, A., SHYAM, P., SASTRY, G., ASKELL, A., ET AL. Language models are few-shot learners. *arXiv preprint arXiv:2005.14165* (2020).
- [8] CAO, L., GAO, L., SONG, J., XU, X., AND SHEN, H. T. Jointly learning attentions with semantic cross-modal correlation for visual question answering. in *Australasian Database Conference* (2017), Springer, pp. 248–260.

- [9] CHEN, K., WANG, J., CHEN, L.-C., GAO, H., XU, W., AND NEVATIA, R. Abc-cnn: An attention based convolutional neural network for visual question answering. *arXiv preprint arXiv:1511.05960* (2015).
- [10] CHEN, Y.-C., LI, L., YU, L., EL KHOLY, A., AHMED, F., GAN, Z., CHENG, Y., AND LIU, J. Uniter: Universal image-text representation learning. in *European Conference on Computer Vision* (2020), Springer, pp. 104–120.
- [11] CHO, K., VAN MERRIËNBOER, B., GULCEHRE, C., BAHDANAU, D., BOUGARES, F., SCHWENK, H., AND BENGIO, Y. Learning phrase representations using rnn encoder-decoder for statistical machine translation. in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (oct 2014), Association for Computational Linguistics, pp. 1724–1734.
- [12] CHOROMANSKI, K., LIKHOSHERSTOV, V., DOHAN, D., SONG, X., GANE, A., SARLÓS, T., HAWKINS, P., DAVIS, J., MOHIUDDIN, A., KAISER, L., BELANGER, D., COLWELL, L. J., AND WELLER, A. Rethinking attention with performers. *ArXiv abs/2009.14794* (2020).
- [13] DENG, J., DONG, W., SOCHER, R., LI, L.-J., LI, K., AND FEI-FEI, L. Imagenet: A large-scale hierarchical image database. in *2009 IEEE conference on computer vision and pattern recognition* (2009), Ieee, pp. 248–255.
- [14] DENKOWSKI, M., AND LAVIE, A. Meteor universal: Language specific translation evaluation for any target language. in *Proceedings of the ninth workshop on statistical machine translation* (2014), pp. 376–380.
- [15] DEVLIN, J., CHANG, M.-W., LEE, K., AND TOUTANOVA, K. Bert: Pre-training of deep bidirectional transformers for language understanding. in *NAACL-HLT* (2019).
- [16] FUKUI, A., PARK, D. H., YANG, D., ROHRBACH, A., DARRELL, T., AND ROHRBACH, M. Multimodal compact bilinear pooling for visual question answering and visual grounding. in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing* (nov 2016), Association for Computational Linguistics, pp. 457–468.
- [17] GAO, H., MAO, J., ZHOU, J., HUANG, Z., WANG, L., AND XU, W. Are you talking to a machine? dataset and methods for multilingual image question. in *Advances in neural information processing systems* (2015), pp. 2296–2304.

- [18] GONG, Y., KE, Q., ISARD, M., AND LAZEBNIK, S. A multi-view embedding space for modeling internet images, tags, and their semantics. *International journal of computer vision* 106, 2 (2014), 210–233.
- [19] GOYAL, Y., KHOT, T., SUMMERS-STAY, D., BATRA, D., AND PARIKH, D. Making the v in vqa matter: Elevating the role of image understanding in visual question answering. in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2017), pp. 6904–6913.
- [20] GURARI, D., LI, Q., STANGL, A. J., GUO, A., LIN, C., GRAUMAN, K., LUO, J., AND BIGHAM, J. P. Vizwiz grand challenge: Answering visual questions from blind people. in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2018), pp. 3608–3617.
- [21] HE, K., ZHANG, X., REN, S., AND SUN, J. Deep residual learning for image recognition. in *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016), pp. 770–778.
- [22] HINTON, G. E., KRIZHEVSKY, A., AND SUTSKEVER, I. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems* 25 (2012), 1106–1114.
- [23] HOCHREITER, S., AND SCHMIDHUBER, J. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [24] JABRI, A., JOULIN, A., AND VAN DER MAATEN, L. Revisiting visual question answering baselines. in *European conference on computer vision* (2016), Springer, pp. 727–739.
- [25] JOHNSON, J., HARIHARAN, B., VAN DER MAATEN, L., FEI-FEI, L., LAWRENCE ZITNICK, C., AND GIRSHICK, R. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2017), pp. 2901–2910.
- [26] KAFLE, K., AND KANAN, C. Answer-type prediction for visual question answering. in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2016), pp. 4976–4984.
- [27] KAFLE, K., AND KANAN, C. An analysis of visual question answering algorithms. in *Proceedings of the IEEE International Conference on Computer Vision* (2017), pp. 1965–1973.

- [28] KAFLE, K., YOUSEFHUSSEN, M., AND KANAN, C. Data augmentation for visual question answering. in *Proceedings of the 10th International Conference on Natural Language Generation* (2017), pp. 198–202.
- [29] KIM, J.-H., LEE, S.-W., KWAK, D., HEO, M.-O., KIM, J., HA, J.-W., AND ZHANG, B.-T. Multi-modal residual learning for visual qa. *Advances in neural information processing systems* 29 (2016), 361–369.
- [30] KITAEV, N., KAISER, L., AND LEVSKAYA, A. Reformer: The efficient transformer. *ArXiv abs/2001.04451* (2020).
- [31] KRISHNA, R., ZHU, Y., GROTH, O., JOHNSON, J., HATA, K., KRAVITZ, J., CHEN, S., KALANTIDIS, Y., LI, L.-J., SHAMMA, D. A., ET AL. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision* 123, 1 (2017), 32–73.
- [32] LAO, M., GUO, Y., WANG, H., AND ZHANG, X. Cross-modal multistep fusion network with co-attention for visual question answering. *IEEE Access* 6 (2018), 31516–31524.
- [33] LI, G., DUAN, N., FANG, Y., GONG, M., JIANG, D., AND ZHOU, M. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. in *AAAI* (2020), pp. 11336–11344.
- [34] LI, X., YIN, X., LI, C., ZHANG, P., HU, X., ZHANG, L., WANG, L., HU, H., DONG, L., WEI, F., ET AL. Oscar: Object-semantics aligned pre-training for vision-language tasks. in *European Conference on Computer Vision* (2020), Springer, pp. 121–137.
- [35] LIN, T.-Y., MAIRE, M., BELONGIE, S., HAYS, J., PERONA, P., RAMANAN, D., DOLLÁR, P., AND ZITNICK, C. L. Microsoft coco: Common objects in context. in *European conference on computer vision* (2014), Springer, pp. 740–755.
- [36] LIN, X., AND PARIKH, D. Leveraging visual question answering for image-caption ranking. in *European Conference on Computer Vision* (2016), Springer, pp. 261–277.
- [37] LIOUTAS, V., PASSALIS, N., AND TEFAS, A. Explicit ensemble attention learning for improving visual question answering. *Pattern Recognition Letters* 111 (2018), 51–57.
- [38] LU, J., BATRA, D., PARIKH, D., AND LEE, S. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. in *Advances in Neural Information Processing Systems* (2019), pp. 13–23.

- [39] LU, J., YANG, J., BATRA, D., AND PARIKH, D. Hierarchical question-image co-attention for visual question answering. in *Advances in neural information processing systems* (2016), pp. 289–297.
- [40] MA, L., LU, Z., AND LI, H. Learning to answer questions from image using convolutional neural network. in *AAAI* (2016).
- [41] MALINOWSKI, M., DOERSCH, C., SANTORO, A., AND BATTAGLIA, P. Learning visual question answering by bootstrapping hard attention. in *Proceedings of the European Conference on Computer Vision (ECCV)* (2018), pp. 3–20.
- [42] MALINOWSKI, M., AND FRITZ, M. A multi-world approach to question answering about real-world scenes based on uncertain input. in *Advances in neural information processing systems* (2014), pp. 1682–1690.
- [43] MALINOWSKI, M., ROHRBACH, M., AND FRITZ, M. Ask your neurons: A deep learning approach to visual question answering. *International Journal of Computer Vision* 125, 1-3 (2017), 110–135.
- [44] MANMADHAN, S., AND KOVOOR, B. C. Visual question answering: a state-of-the-art review. *Artificial Intelligence Review* (2020), 1–41.
- [45] MIKOLOV, T., CHEN, K., CORRADO, G. S., AND DEAN, J. Efficient estimation of word representations in vector space. *CoRR abs/1301.3781* (2013).
- [46] NOH, H., HONGSUCK SEO, P., AND HAN, B. Image question answering using convolutional neural network with dynamic parameter prediction. in *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016), pp. 30–38.
- [47] PAPINENI, K., ROUKOS, S., WARD, T., AND ZHU, W.-J. 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.
- [48] PENG, L., YANG, Y., BIN, Y., XIE, N., SHEN, F., JI, Y., AND XU, X. Word-to-region attention network for visual question answering. *Multimedia Tools and Applications* 78, 3 (2019), 3843–3858.
- [49] PENNINGTON, J., SOCHER, R., AND MANNING, C. D. Glove: Global vectors for word representation. in *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)* (2014), pp. 1532–1543.

- [50] RADFORD, A., WU, J., CHILD, R., LUAN, D., AMODEI, D., AND SUTSKEVER, I. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.
- [51] REN, M., KIROS, R., AND ZEMEL, R. Exploring models and data for image question answering. in *Advances in neural information processing systems* (2015), pp. 2953–2961.
- [52] REN, M., KIROS, R., AND ZEMEL, R. Image question answering: A visual semantic embedding model and a new dataset. *Proc. Advances in Neural Inf. Process. Syst* 1, 2 (2015), 5.
- [53] REN, S., HE, K., GIRSHICK, R., AND SUN, J. Faster r-cnn: Towards real-time object detection with region proposal networks. in *Advances in neural information processing systems* (2015), pp. 91–99.
- [54] RUWA, N., MAO, Q., WANG, L., AND DONG, M. Affective visual question answering network. in *2018 IEEE conference on multimedia information processing and retrieval (MIPR)* (2018), IEEE, pp. 170–173.
- [55] SHAH, S., MISHRA, A., YADATI, N., AND TALUKDAR, P. P. Kvqa: Knowledge-aware visual question answering. in *Proceedings of the AAAI Conference on Artificial Intelligence* (2019), volume 33, pp. 8876–8884.
- [56] SHI, Y., FURLANELLO, T., ZHA, S., AND ANANDKUMAR, A. Question type guided attention in visual question answering. in *Proceedings of the European Conference on Computer Vision (ECCV)* (2018), pp. 151–166.
- [57] SHIH, K. J., SINGH, S., AND HOIEM, D. Where to look: Focus regions for visual question answering. in *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016), pp. 4613–4621.
- [58] SHRESTHA, R., KAFLE, K., AND KANAN, C. Answer them all! toward universal visual question answering models. in *Proceedings of the IEEE conference on computer vision and pattern recognition* (2019), pp. 10472–10481.
- [59] SILBERMAN, N., HOIEM, D., KOHLI, P., AND FERGUS, R. Indoor segmentation and support inference from rgb-d images. in *European conference on computer vision* (2012), Springer, pp. 746–760.
- [60] SIMONYAN, K., AND ZISSERMAN, A. Very deep convolutional networks for large-scale image recognition. *CoRR abs/1409.1556* (2015).

- [61] SU, W., ZHU, X., CAO, Y., LI, B., LU, L., WEI, F., AND DAI, J. VI-bert: Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530* (2019).
- [62] SUN, C., MYERS, A., VONDRICK, C., MURPHY, K., AND SCHMID, C. Videobert: A joint model for video and language representation learning. in *Proceedings of the IEEE International Conference on Computer Vision* (2019), pp. 7464–7473.
- [63] SZEGEDY, C., LIU, W., JIA, Y., SERMANET, P., REED, S., ANGUELOV, D., ERHAN, D., VANHOUCKE, V., AND RABINOVICH, A. Going deeper with convolutions. in *Proceedings of the IEEE conference on computer vision and pattern recognition* (2015), pp. 1–9.
- [64] TALAFHA, B., AND AL-AYYOUB, M. Just at vqa-med: A vgg-seq2seq model. in *CLEF (Working Notes)* (2018).
- [65] TAN, H. H., AND BANSAL, M. Lxmert: Learning cross-modality encoder representations from transformers. in *EMNLP/IJCNLP* (2019).
- [66] TANG, R., MA, C., ZHANG, W. E., WU, Q., AND YANG, X. Semantic equivalent adversarial data augmentation for visual question answering. in *European Conference on Computer Vision* (2020), Springer, pp. 437–453.
- [67] TENEY, D., AND VAN DEN HENGEL, A. Visual question answering as a meta learning task. in *Proceedings of the European Conference on Computer Vision (ECCV)* (2018), pp. 219–235.
- [68] TOMMASI, T., MALLIA, A., PLUMMER, B., LAZEBNIK, S., BERG, A. C., AND BERG, T. L. Combining multiple cues for visual madlibs question answering. *International Journal of Computer Vision* 127, 1 (2019), 38–60.
- [69] TOOR, A. S., WECHSLER, H., AND NAPPI, M. Question action relevance and editing for visual question answering. *Multimedia Tools and Applications* 78, 3 (2019), 2921–2935.
- [70] WANG, P., WU, Q., SHEN, C., DICK, A., AND VAN DEN HENGEL, A. Explicit knowledge-based reasoning for visual question answering. in *Proceedings of the 26th International Joint Conference on Artificial Intelligence* (2017), IJCAI’17, AAAI Press, p. 1290–1296.
- [71] WU, Q., SHEN, C., WANG, P., DICK, A., AND VAN DEN HENGEL, A. Image captioning and visual question answering based on attributes and external knowledge. *IEEE transactions on pattern analysis and machine intelligence* 40, 6 (2017), 1367–1381.

- [72] WU, Q., TENNEY, D., WANG, P., SHEN, C., DICK, A., AND VAN DEN HENGEL, A. Visual question answering: A survey of methods and datasets. *Computer Vision and Image Understanding* 163 (2017), 21–40.
- [73] WU, Z., AND PALMER, M. Verbs semantics and lexical selection. in *Proceedings of the 32nd Annual Meeting on Association for Computational Linguistics* (USA, 1994), ACL '94, Association for Computational Linguistics, p. 133–138.
- [74] XU, H., AND SAENKO, K. Ask, attend and answer: Exploring question-guided spatial attention for visual question answering. in *European Conference on Computer Vision* (2016), Springer, pp. 451–466.
- [75] YANG, Z., HE, X., GAO, J., DENG, L., AND SMOLA, A. Stacked attention networks for image question answering. in *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016), pp. 21–29.
- [76] YU, L., PARK, E., BERG, A. C., AND BERG, T. L. Visual madlibs: Fill in the blank description generation and question answering. in *Proceedings of the IEEE international conference on computer vision* (2015), pp. 2461–2469.
- [77] YU, Z., YU, J., CUI, Y., TAO, D., AND TIAN, Q. Deep modular co-attention networks for visual question answering. in *Proceedings of the IEEE conference on computer vision and pattern recognition* (2019), pp. 6281–6290.
- [78] YU, Z., YU, J., XIANG, C., FAN, J., AND TAO, D. Beyond bilinear: Generalized multimodal factorized high-order pooling for visual question answering. *IEEE transactions on neural networks and learning systems* 29, 12 (2018), 5947–5959.
- [79] ZHOU, L., PALANGI, H., ZHANG, L., HU, H., CORSO, J. J., AND GAO, J. Unified vision-language pre-training for image captioning and vqa. in *AAAI* (2020), pp. 13041–13049.
- [80] ZHU, Y., GROTH, O., BERNSTEIN, M., AND FEI-FEI, L. Visual7w: Grounded question answering in images. in *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016), pp. 4995–5004.



# واژه‌نامه فارسی به انگلیسی

data augmentation . . . . .	افزایش داده
bias . . . . .	بایاس
text-to-image retrieval . . . . .	بازیابی متن به تصویر
global vector . . . . .	بردار سراسری
computer vision . . . . .	بینایی ماشین
natural language processing . . . . .	پردازش زبان طبیعی
sentiment analysis . . . . .	تحلیل احساسات
object detection . . . . .	تشخیص اشیا
activity detection . . . . .	تشخیص فعالیت
machine translation . . . . .	ترجمه ماشینی
image captioning . . . . .	توصیف تصویر
object segmentation . . . . .	تقسیم‌بندی اشیا
face recognition . . . . .	تشخیص چهره
word embedding . . . . .	تعبیه کلمات
hard attention . . . . .	توجه سخت
soft attention . . . . .	توجه نرم
token . . . . .	توکن
single-stream . . . . .	تک جریان
element-wise addition . . . . .	جمع متناظر
rotation . . . . .	چرخش
annotation . . . . .	حاشیه‌نویسی
text summarization . . . . .	خلاصه‌سازی متون
voice assistants . . . . .	دستیاران صوتی
two-stream . . . . .	دو جریان
recommender systems . . . . .	سیستم‌های توصیه‌گر
artificial neural networks . . . . .	شبکه‌های عصبی مصنوعی

convolutional neural networks	شبکه‌های عصبی پیچشی
fully connected networks	شبکه‌های کاملاً متصل
recurrent neural networks	شبکه‌های عصبی بازگشتی
element-wise multiplication	ضرب متناظر
object classification	طبقه‌بندی اشیا
attribute classification	طبقه‌بندی صفات
scene classification	طبقه‌بندی صحنه
text classification	طبقه‌بندی متون
conversational agents	عامل‌های گفتگو
spell correction	غلط‌یابی متون
vanishing gradient	محوشدگی گرادیان
encoder-decoder architecture	معماری رمزگذار-رمزگشا
attention mechanism	مکانیزم توجه
adversarial examples	نمونه‌های خصمانه
deep learning	یادگیری عمیق
machine learning	یادگیری ماشین

# واژه‌نامه انگلیسی به فارسی

activity detection	تشخیص فعالیت
annotation	حاشیه‌نویسی
attention mechanism	مکانیزم توجه
attribute classification	طبقه‌بندی صفات
artificial neural networks	شبکه‌های عصبی مصنوعی
adversarial examples	نمونه‌های خصمانه
bias	بایاس
computer vision	بینایی ماشین
conversational agents	عامل‌های گفتگو
convolutional neural networks	شبکه‌های عصبی پیچشی
data augmentation	افزایش داده
deep learning	یادگیری عمیق
encoder-decoder architecture	معماری رمزگذار-رمزگشا
element-wise multiplication	ضرب متناظر
element-wise addition	جمع متناظر
face recognition	تشخیص چهره
fully connected networks	شبکه‌های کاملاً متصل
global vector	بردار سراسری
hard attention	توجه سخت
image captioning	توصیف تصویر
machine translation	ترجمه ماشینی
machine learning	یادگیری ماشین
natural language processing	پردازش زبان طبیعی
object detection	تشخیص اشیا
object segmentation	تقسیم‌بندی اشیا
object classification	طبقه‌بندی اشیا

recommender systems	سیستم‌های توصیه‌گر
recurrent neural networks	شبکه‌های عصبی بازگشتی
rotation	چرخش
sentiment analysis	تحلیل احساسات
soft attention	توجه نرم
single-stream	تک جریان
scene classification	طبقه‌بندی صحنه
spell correction	غلط‌یابی متون
text-to-image retrieval	بازیابی متن به تصویر
text summarization	خلاصه‌سازی متون
text classification	طبقه‌بندی متون
two-stream	دو جریان
token	توکن
voice assistants	دستیاران صوتی
vanishing gradient	محوشدگی گرادیان
word embedding	تعبیه کلمات

**Abstract:**

Visual Question Answering(VQA) is a challenging task that has been introduced in recent years and has received increasing attention from both the computer vision and the natural language processing communities. Visual Question Answering aims to answer the questions about given images. A VQA system tries to find the correct answer to questions using visual elements of the image and inference gathered from textual questions. In the first chapter of this review, we present the Visual Question Answering task, applications, and challenges. After defining some concepts in the second chapter, we discuss various datasets for VQA, methods, and evaluation metrics in chapter 3. Due to the success of deep learning and pre-trained models, we classify VQA methods into two general approaches: deep learning and pre-trained models. In the last chapter, after concluding on the different aspects of VQA, we provide some directions for future work.

**Keywords:** Visual Question Answering, Natural Language Processing, Computer Vision, Deep Learning, pretrained models



**Iran University of Science and Technology  
Computer Engineering Department**

# **Visual Question Answering**

**A Thesis Submitted in Partial Fulfillment of the Requirement for the Degree  
of Master of Science in Computer Engineering**

**By:**

Maryam Sadat Hashemi

**Supervisor:**

**Dr. Sayyed Sauleh Eetemadi**

**December 2020**