

بهبود پاسخ دهی به متون درون عکس با تصحیح توکن‌های استخراج شده از تصویر

کبری فرشیدی*، حسن ختن لو، محرم منصوری زاده

گروه مهندسی کامپیوتر، دانشکده مهندسی، دانشگاه بوعلی سینا، همدان

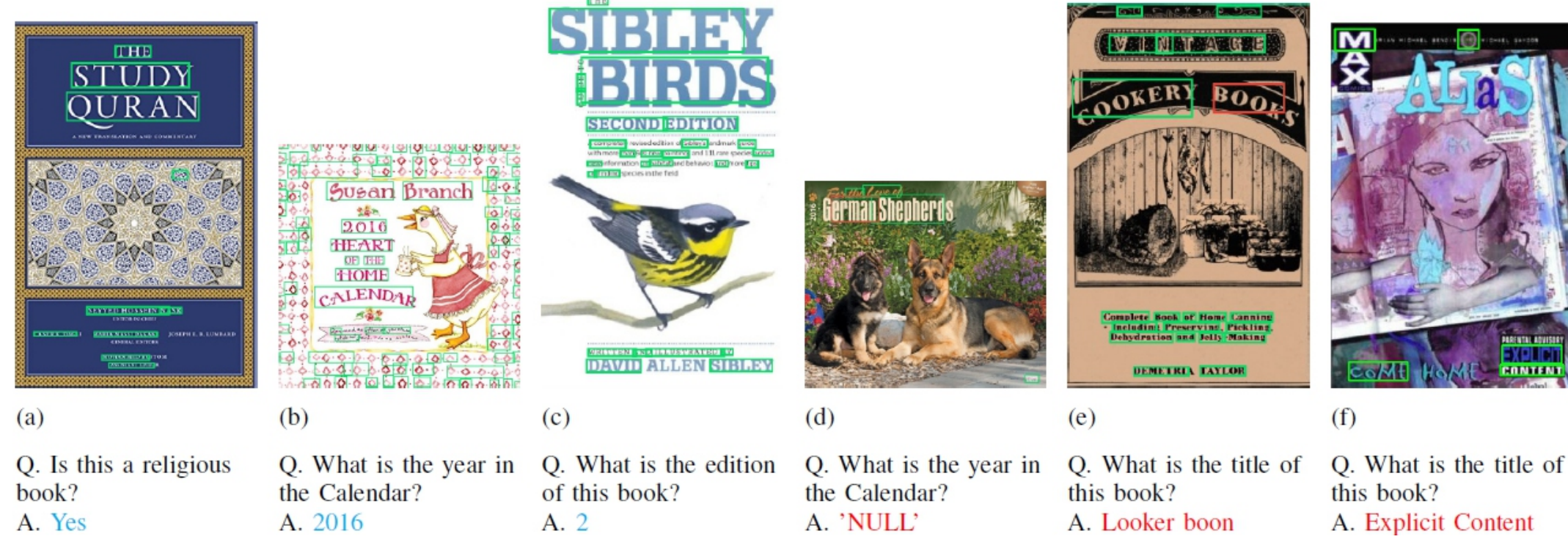
k.farshidi@eng.basu.ac.ir

خلاصه:

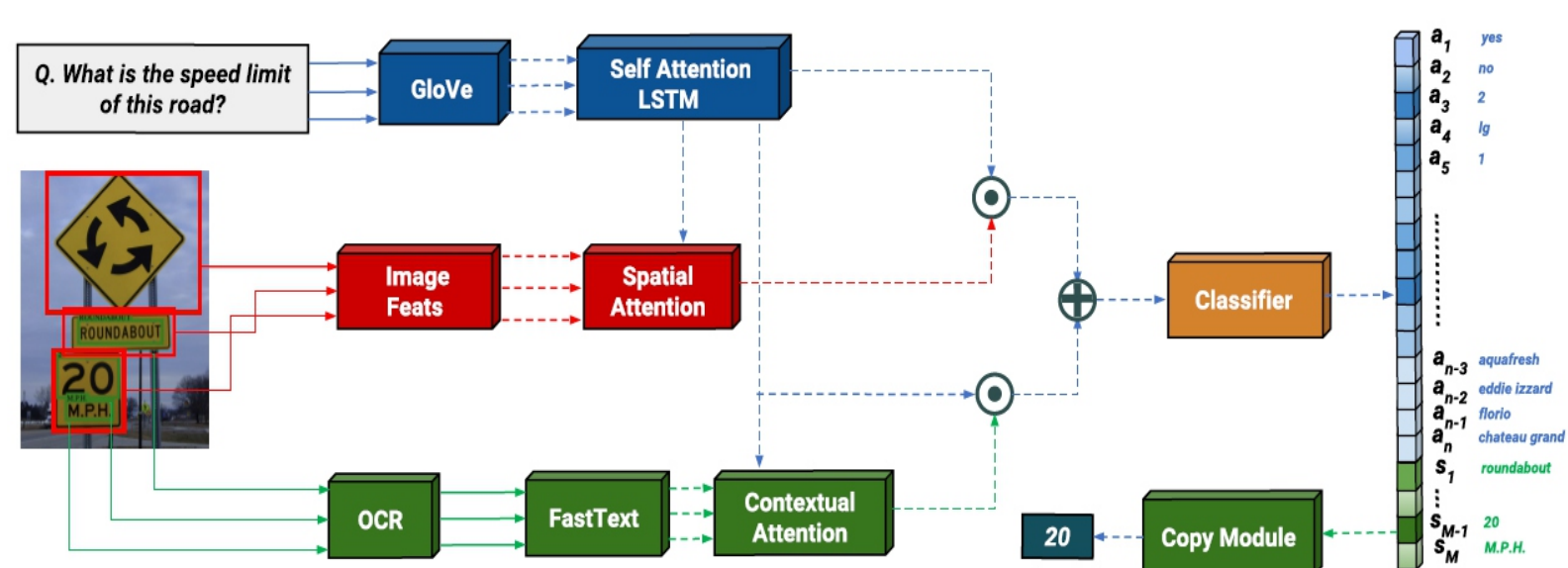
یکی از انواع مسائل چند وجهی، مسئله‌ی پاسخ‌گویی به متون درون تصویر است (شکل ۱). هدف، کشف ارتباط بین تصویر (بخصوص متون و علائم درون تصویر) و سوال متنی پرسیده شده، درباره‌ی آن تصویر می‌باشد. در نتیجه‌ی این پژوهش ساخت دستیار هوشمند برای معلولان، نابینایان و کودکان خردسال است حتی از ثمره‌ی این پژوهش میتوان در ماشین‌های تسلا بهره برد. پژوهش‌گران توسط مدل‌های یادگیری عمیق، معماری مبدل‌ها، مکانیزم توجه و ... توانسته‌اند به نتایج خوبی دست پیدا کنند. اما هنوز این حوزه بصورت عام و گسترده بکار برده نشده است. چرا که دقت به حد قابل قبولی نرسیده است. یکی از دلایل دقت پایین، ورود توکن‌های استخراج شده‌ی اشتباه در مدل میباشد. چرا که موتور استخراج‌کننده‌ی توکن هر چقدر هم قوی بوده باشد، باز هم دارای خطا در تولید توکن‌ها است. در این پژوهش قصد داریم با به کار گیری پیش پردازشها و تصحیح توکن‌های ورودی این دقت را بالا ببریم.

کلمات کلیدی:

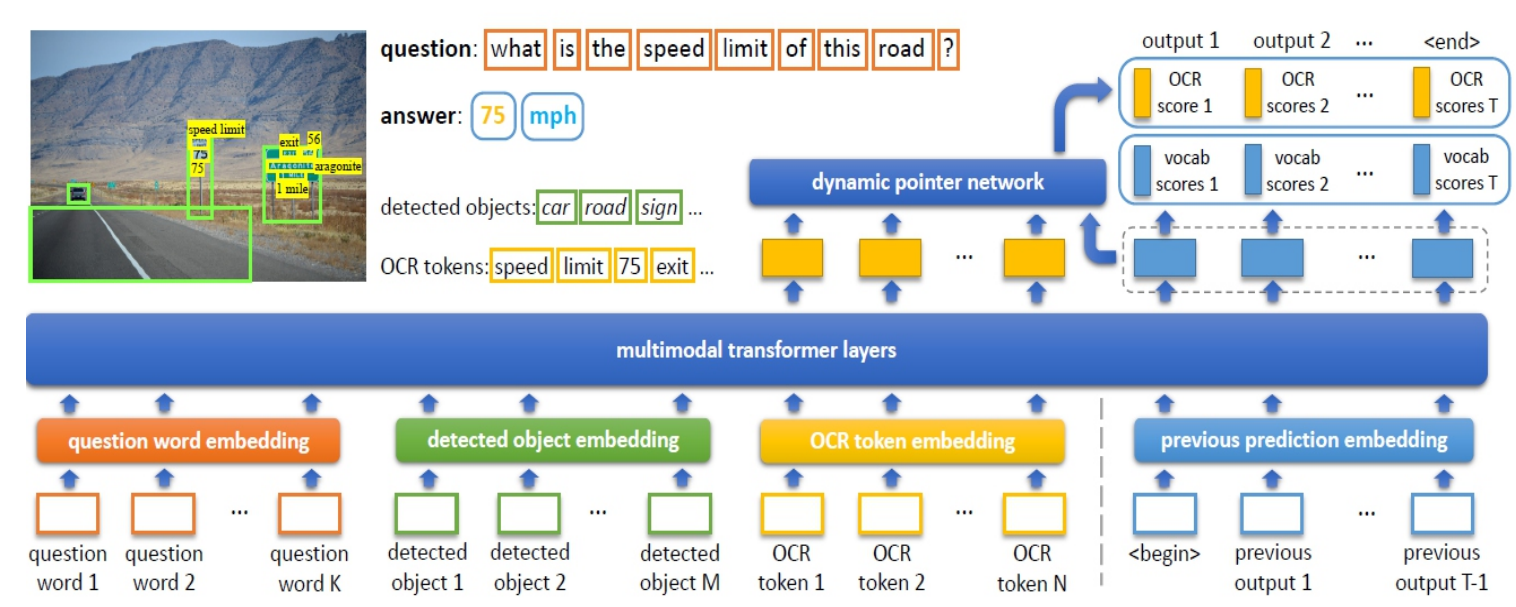
پرسش از متون درون تصویر، مبدل، مکانیزم توجه، مدل از پیش آموزش دیده



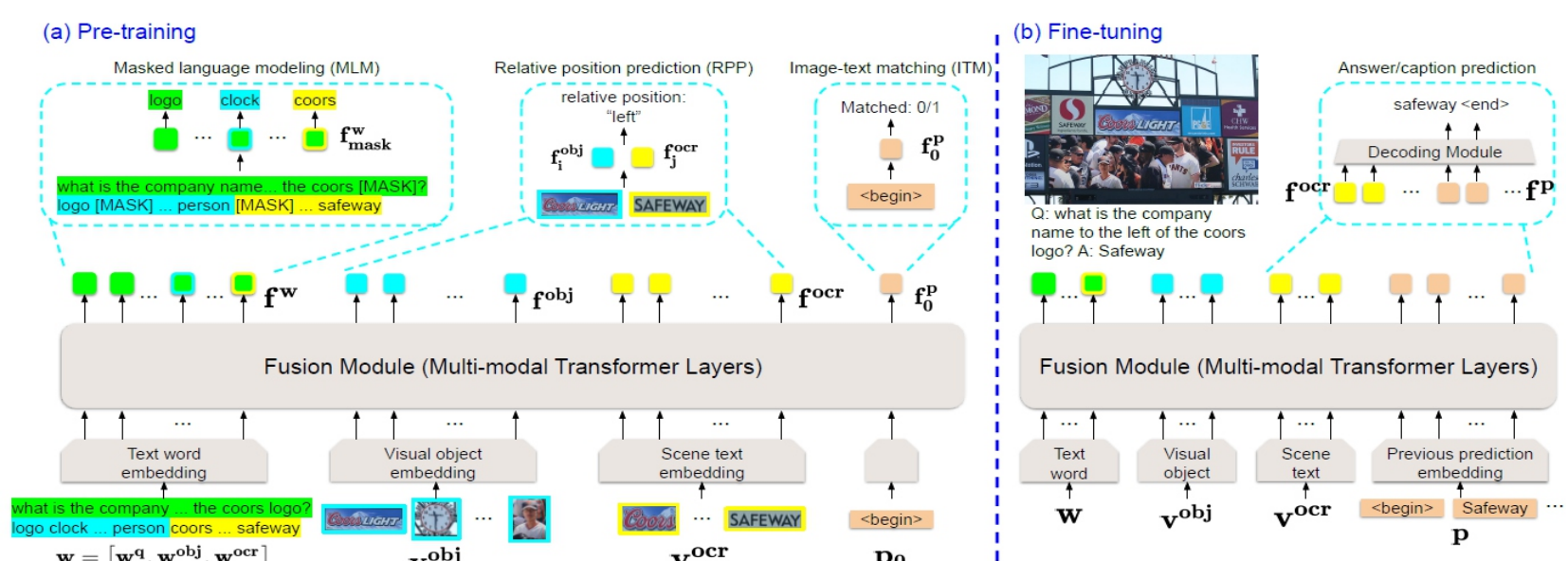
شکل ۱:



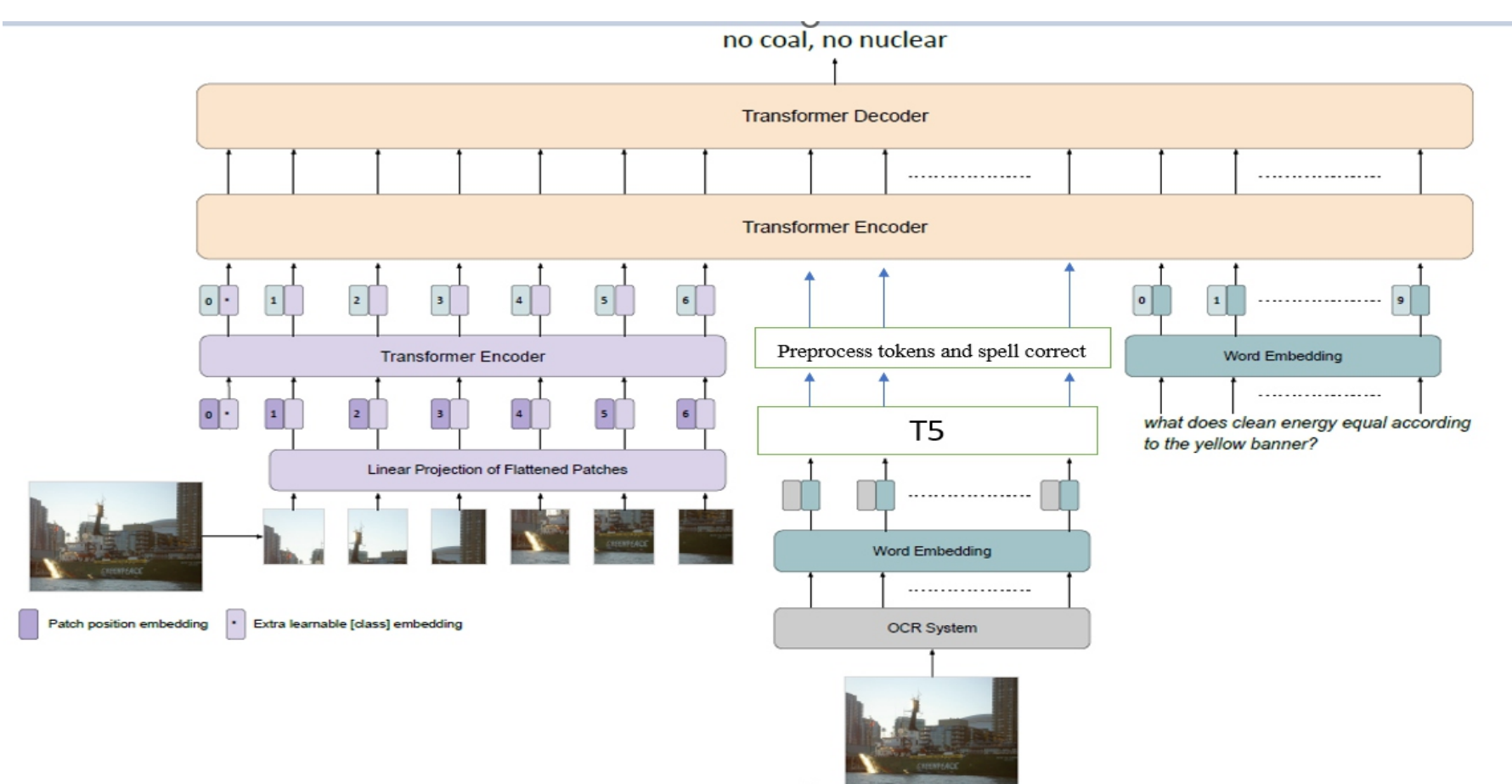
شکل ۲:



شکل ۳:



شکل ۴:



شکل ۵:

مواد و روشها:

همانطور که گفته شد ما از دسته سوم پژوهشگران بهره میبریم و دقت آن را بالاتر می‌بریم. همانند شکل روبرو (شکل ۵) به این صورت عمل می‌کنیم که قبل از اینکه ویژگی‌های استخراج شده وارد معماری مبدل شود و همه‌ی وجه‌ها با هم ادغام شوند، ویژگی‌های توکن‌ها را بهبود می‌دهیم چرا که در صورتی که توکن‌های اشتباه وارد معماری مبدل شود قطعاً جواب‌های اشتباهی نیز استخراج میشود. و استدلال‌های اشتباهی از متون درون تصویر میشود. مدل‌های پیش آموزش دیده‌ی متعددی به جهت بررسی صحت توکن‌ها و جملات بوجود آمده است که میتوان با استفاده از آنها توکن‌ها را بهبود داد و از طرفی با یک پیش پردازش متون استخراج شده از تصویر جملات بهتری برای استدلال بوجود آورد.

بحث و نتیجه‌گیری:

چالش اساسی پژوهش‌گران بیشتر آماده‌سازی مجموعه داده‌های مرتبط با این پژوهش، مجموعه داده ساده‌تر و در آن واحد دارای اطلاعات غنی‌تر، طراحی معماری‌های جدیدتر با توجه به پیشرفت معماری‌های شبکه‌های عمیق برای بالابردن عملکرد این حوزه میباشد. ما در معماری خود سعی کرده ایم با بهره‌گیری بهتر از توکن‌های استخراج شده از طریق موتور جستجوگر، دقت این حوزه را بالا ببریم.

منابع و مراجع:

- [1] amanpreet S., Vivek N., Meet S., Yu J., Xinlei C., Dhruv B., Devi P. and Marcus R., Towards VQA model that can read. IEEE 2020
- [2] Ronghang Hu., Amanpreet S., Trevor D., Iterative answer prediction with pointer-argued multimodal transformer or TVQA. IEEE 2020
- [3] Chenyu G., Qi Z., Peng Wang, Hui L., Yuliang L., Anton V. . Structured multimodal attention for TVQA. arXiv.2020
- [4] Zhengyuan Y., Yijuan L., Jianfeng W., Xi Y., Dinei F., Lijuan W., Cha Z., Lei Z., TAP: Text-aware pretraining for TVQA and Text caption. arXiv 2020
- [5] Furkan A., Litman R., Xie Y., appalaraju S., . LaTr: Layout-aware Tranformer for Scene-Text VQA. arXiv.2021

Visual Question Answering=VQA
Text Based Visual Question Answering=TVQA
Optical Character Recognition=OCR

لیست اختصارات: