



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

دانشکده مهندسی کامپیوتر و فناوری اطلاعات

**پرسش و پاسخ بصری با استفاده از شبکه های عصبی
کانولوشنی و بازگشتی عمیق**

علی غلامی

۱۶ اردیبهشت ۱۳۹۷

۱ مقدمه

درك انسان از محیط اطراف خود، از ابتدای کودکی شکل می گیرد. اگر از یک کودک، درباره جهانی که برای او قابل رویت است، سوالی پرسیده شود، بلافاصله درباره این جهان جملاتی توصیفی با دقت بسیار بالا ارائه می کند. انتقال این قدرت به ماشین ها از این جهات مختلفی حائز اهمیت می باشد. یکی از جهات اصلی آن، کاربرد های بی شماری است که می توان از این قابلیت در امور مختلفی استفاده کرد. دوربین های نظارتی، خودروهای خودران و ... نمونه هایی از این کاربرد ها هستند. درک تصاویر به همراه درک سوالاتی که مرتبط با آن تصاویر پرسیده می شود، می تواند یکی از اساسی ترین قابلیت های یک سامانه مدیریت هوشمند تصاویر، خودروی خودران و یا یک سیستم بازیابی تصویر هوشمند باشد. برای این منظور نیاز است که ابتدا صحنه به نمایش درآمده توسط ماشین هضم گردد. بدین معنی که بدون درک درست از آنچه در تصویر موجود است، ساخت چنین سامانه ای امکان پذیر نخواهد بود

۲ مرور سوابق پیشین

روش هایی که جهت حل مسائل پرسش و پاسخ بصري ارائه شده اند از تنوع بالایی برخوردار هستند. اگرچه، می توان اغلب این روش ها را در راستای حل چهار چالش مهم زیر در نظر گرفت:

۱. استخراج ویژگی از تصاویر
۲. استخراج ویژگی از متن و سوالات پرسیده شده
۳. ترکیب بردار های ویژگی مستخرج از تصویر و متن
۴. چالش تولید پاسخ متناسب با ویژگی ها ترکیب شده

با تحول عظیمی که در سال ۲۰۱۲ با ارائه مدل کانولوشنی عمیق جهت استخراج ویژگی ارائه شد و نیز افزایش قدرت پردازشی ماشین ها، توجه به سمت شبکه های کانولوشنی جهت استخراج ویژگی از تصاویر بیشتر شد. با پیشرفت این مدل ها و ترکیب آنها با مدل های پردازش زبان طبیعی، تولید شرح بر تصاویر نیز مورد توجه قرار گرفت. از سال ۲۰۱۵ تا به

اکنون، یکی از جالب ترین موضوعاتی که توجه پژوهشگران را به سمت خود جلب کرده است، موضوع پرسش و پاسخ بصري می باشد. این موضوع ارتباط تنگاتنگی با موضوع تولید شرح بر تصاویر دارد. به همین منظور، بسیاری از تکنیک های رایج در بحث شرح بر تصاویر، در موضوع پرسش و پاسخ بصري نیز مورد توجه قرار گرفته است.

۳ طرح پیشنهادی

سیستم پرسش و پاسخ بصري مطرح در این گزارش، به زبان پایتون و با استفاده از چارچوب کاری تنسورفلو پیاده سازی خواهد شد. هسته این چارچوب کاری به زبان سی پلاس پلاس و با استفاده از پلتفرم توسعه موازی کودا پیاده سازی شده است. این چارچوب کاری توسط تیم گوگل برین در حال توسعه بوده و پشتیبانی می شود. ایده ی اصلی مطرح در این چارچوب کاری، بیان محاسبات در قالب گراف است. هر گره ی این گراف، یک واحد محاسباتی را مشخص می کند. با این رویکرد می توان شبکه های پیچیده را به راحتی پیاده سازی نمود.

۴ محصولات طرح

خروجی های این طرح شامل موارد زیر خواهند بود:

- بازیابی تصاویر در شبکه های اجتماعی با استفاده از جستجوی متن
- توصیف کننده جهان اطراف جهت استفاده نا بینایان و مذاکره دو طرفه
- ربات های امداد گر در عملیات هایی که با انسان تبادل اطلاعات انجام می دهند
- ربات های مذاکره کننده

۵ مراحل انجام

۱. شناسایی و تهیه منابع:

- پیدا کردن مقالات مربوط به پرسش و پاسخ بصری
- یافتن توضیحات مربوط به طرح های از پیش پیاده سازی شده

۲. تنظیم ساختار:

- تهیه ی بخشهای مقدمه، محتوای اصلی، چیکده، نتیجه گیری و منابع
- مشخص کردن ترتیب مباحث محتوا
- تهیه ی فهرست مباحث اصلی و فرعی یافتن معماری شبکه های کانولوشنی

۳. مطالعه و یادداشت برداری:

- پیدا کردن مقالات مربوط به پرسش و پاسخ بصری
- یافتن توضیحات مربوط به طرح های از پیش پیاده سازی شده

۴. اجرای بخش عملی:

- پیاده سازی شبکه های کانولوشنی عمیق
- پیاده سازی شبکه بازگشتی عمیق
- آموزش شبکه
- آزمایش و ارزیابی

۵. تهیه گزارش نهایی:

- مستند سازی و نوشتن گزارش نهایی
- آمادگی جهت ارائه ی شفاهی

۶ زمان بندی پروژه

زمان بندی پروژه بر اساس “پیشنهاد زمان بندی پروژه کارشناسی” در کتاب انجام شده است.

