

طراحی و پیاده‌سازی کپچای فارسی با استفاده از ویژگی‌های دستور

زبان فارسی

فرزین یغمایی^۱، حمید نصیری^۲

^۱دانشگاه سمنان، هیات علمی دانشکده مهندسی برق و کامپیوتر ، f_yaghmaee@semnan.ac.ir

^۲دانشگاه سمنان، دانشکده مهندسی برق و کامپیوتر ، nasiri.hamid@gmail.com

چکیده

امروزه با گسترش اینترنت استفاده از خدمات تحت وب نیز گسترش یافته‌است. یکی از مشکلاتی که این خدمات با آن مواجه هستند، برنامه‌های کامپیوتری موسوم به ربات‌های وب‌گرد می‌باشند. برای جلوگیری از حملات ربات‌های وب‌گرد و همچنین کاهش هدر رفتن منابع توسط آن‌ها از روشی به نام کپچا^۱ استفاده می‌شود که به سیستم‌های کامپیوتری کمک می‌کند تا به تفکیک کاربران انسانی از نرم‌افزارهای رایانه‌ای بپردازند. در سال‌های اخیر با توجه به گسترش استفاده از زبان فارسی در وب، طراحی کپچاهای فارسی امری ضروری به حساب می‌آید. در این مقاله سعی شده‌است تا با استفاده از دستور زبان فارسی، روش جدیدی در طراحی کپچای فارسی با امنیت و کارایی بالا ارائه شود. در روش پیشنهاد شده جمله‌ای در تصویر به کاربر نمایش داده می‌شود که از نظر دستوری (فعل یا فاعل) اشکال دارد که کاربر باید فعل یا فاعل جمله را اصلاح نماید. به منظور افزایش سطح امنیت کپچای ارائه شده، اندازه قلم و محل قرارگیری جمله در تصویر به صورت تصادفی انتخاب می‌شوند و به تصویر نهایی تعدادی خط به عنوان نویز اضافه می‌شود. نتایج بدست آمده نشان می‌دهند که کاربر فارسی زبان به راحتی قادر به اصلاح جمله از لحاظ دستوری می‌باشد. در حالی که این کار برای ماشین به سختی امکان پذیر است.

واژه‌های کلیدی

کپچا، کپچای فارسی، نرم‌افزارهای تشخیص متن، دستور زبان فارسی

۱- مقدمه

امروزه با گسترش اینترنت در سراسر جهان، استفاده از خدمات تحت وب توسط کاربران نیز گسترش یافته‌است. برای استفاده از این خدمات کاربر باید اطلاعات خود را از طریق فرم‌های موجود در صفحات وب ارسال نماید. اما مسئله به همین سادگی نیست و مدیران وب‌سایت‌ها با چالش‌های مختلفی در این زمینه روبرو هستند که یکی از آن‌ها ربات‌های وب‌گرد می‌باشند.

ربات‌های وب‌گرد برنامه‌های کامپیوتری هستند که با پرکردن فرم‌ها به صورت خودکار و با سرعت زیاد مشکلاتی را برای وب‌سایت‌ها ایجاد می‌کنند. این ربات‌ها آنقدر گسترش یافته‌اند که براساس گزارش شرکت Incapsula حدود ۶۱٪ ترافیک مصرفی اینترنت مربوط به آن‌ها است که این میزان نسبت به سال ۲۰۱۲، ۲۱٪ افزایش داشته‌است [۴].

برای جلوگیری از این مشکلات معمولاً از کپچا استفاده می‌شود. کپچا هر نوع تستی است که به صورت خودکار تولید شده و می‌تواند بین کاربران انسانی و ربات‌های کامپیوتری تمایز ایجاد کند. به صورتی که اغلب انسان‌ها می‌توانند از این تست عبور کنند ولی برنامه‌های کامپیوتری فعلی قادر به عبور از آن نیستند. کپچا از این لحاظ که بین انسان و کامپیوتر تمایز قائل می‌شود، شبیه آزمون تورینگ است. اما تفاوت آن با آزمون تورینگ در این

است که در اینجا داور کامپیوتر است [۵]. ساخت یک کپچای جدید وضعیتی برد-برد است. به این معنی که اگر شکسته نشود راهی جدید برای تمایز میان انسان و کامپیوتر پیدا شده و اگر شکسته شود مسئله‌ای در حوزه هوش مصنوعی حل شده‌است [۶].

در سال‌های اخیر برخی نرم افزارهای تشخیص متن^۲ توانسته‌اند نمونه‌هایی از کپچا را بشکنند. برخی شرکت‌های فعال در حوزه هوش مصنوعی نیز ادعا کرده‌اند که توانسته‌اند به طور کلی کپچا را از میان بردارند. اما این ادعاها تازگی ندارد. به گفته پروفیسور لوییس وون آهن (شخصی که برای اولین بار واژه CAPTCHA را معرفی کرد [۷])، از سال ۲۰۰۳ هر ۶ ماه یکبار شخصی ادعا می‌کند که کامپیوترها می‌توانند کپچا را به طور کامل بشکنند [۸]. اما همچنان کپچا از محبوبیت زیادی در دنیای اینترنت برخوردار است. به طوری که روزانه ۲۸۰ میلیون کپچا در سراسر جهان تست می‌شود [۹].

برخلاف گسترش و تکامل نرم افزارهای تشخیص متن در زبان لاتین تاکنون برای زبان فارسی نرم افزار تشخیص متن قدرتمندی تولید نشده است که یکی از دلایل آن پیچیدگی بالا و مشکل بودن ساختار و نوشتار زبان فارسی در مقایسه با زبان لاتین می‌باشد [۱]. از این رو تشخیص کپچای فارسی در مقایسه با کپچای انگلیسی برای نرم افزارهای تشخیص متن مشکل‌تر بوده و قابلیت اطمینان آن در مقابل حمله ربات‌ها بیشتر

سایت PayPal برای ساخت کپچا مطابق شکل ۳ از ترکیب تصادفی حروف و اعداد که با فاصله زیاد از هم قرار گرفته بودند، استفاده می‌کرد. بر روی این کاراکترها خطوط افقی و عمودی نمایش داده می‌شد تا تشخیص آن برای کامپیوتر مشکل شود ولی به دلیل فاصله زیاد کاراکترها از هم به راحتی شکسته شد [۱۳].



شکل ۳: نمونه کپچای سایت PayPal

شرکت Google از کپچایی تحت عنوان reCAPTCHA برای دیجیتالی کردن کتاب‌ها و روزنامه‌های قدیمی استفاده می‌نماید. در این کپچا دو لغت به کاربر نمایش داده می‌شود که باید به هردوی آن‌ها پاسخ دهد. یکی از آن‌ها از مجموعه لغاتی است که در دیجیتالی کردن کتاب‌ها توسط OCR تشخیص داده نشده است و دیگری لغت کنترلی است که کامپیوتر پاسخ آن را می‌داند. در صورتی که کاربر به لغت کنترلی پاسخ درست داده باشد، پاسخ او برای لغت دیگر به عنوان یکی از پاسخ‌های احتمالی در نظر گرفته می‌شود. اگر تعداد زیادی از کاربران به یک لغت پاسخ یکسانی بدهند، آن پاسخ برای دیجیتالی کردن کتاب‌ها در اختیار نرم‌افزار OCR قرار می‌گیرد [۱۴]. نمونه‌ای از reCAPTCHA در شکل ۴ مشاهده می‌شود.



شکل ۴: نمونه reCAPTCHA شرکت Google [۱۵]

علاوه بر موارد ذکر شده در بالا کارهای دیگری نیز در زمینه کپچاهای مبتنی بر کاراکتر انجام شده است که از آن جمله می‌توان به Pessimial [۱۶] و Baffle Text Method [۱۷] اشاره کرد.

۲-۲- کپچاهای مبتنی بر تصویر

در این نوع کپچا تصویر یا تصاویری به کاربر نمایش داده می‌شود که کاربر باید بتواند محتوای تصویر یا تصاویر را تشخیص بدهد. کپچاهای مبتنی بر تصویر معمولاً کاربر پسندتر از کپچاهای مبتنی بر کاراکتر هستند [۱۸] و تشخیص آن‌ها برای ربات‌های کامپیوتری در اکثر مواقع سخت‌تر از تشخیص کپچاهای مبتنی بر کاراکتر است [۱۲]. از کارهای انجام شده در این زمینه می‌توان به ASIRRA^۳ اشاره کرد. ASIRRA کپچای شرکت Microsoft است که در آن ۱۲ تصویر گربه و سگ مطابق شکل ۵ به کاربر نمایش داده می‌شود که کاربر باید تمامی تصاویر گربه‌ها را انتخاب نماید [۱۹].

است. به همین علت در این مقاله سعی شده با بهره‌گیری از ویژگی‌های دستوری زبان فارسی، کپچای فارسی با قابلیت اطمینان بالا و استفاده آسان‌تر برای کاربران فارسی زبان ارائه شود.

ادامه مقاله بدین صورت سازماندهی شده است. در بخش دوم مروری بر کارهای دیگران در زمینه کپچا خواهیم داشت. در بخش سوم به ویژگی‌های منحصر به فرد زبان فارسی، ساختار جملات فارسی و کارهای انجام شده بر روی کپچای فارسی می‌پردازیم. در بخش چهارم روش پیشنهادی در این مقاله را شرح می‌دهیم. در بخش پنجم نتایج حاصل از روش پیشنهادی را ارائه می‌دهیم و در نهایت در بخش ششم به نتیجه‌گیری و جمع‌بندی طرح ارائه شده می‌پردازیم.

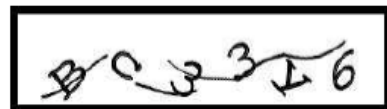
۲- کارهای دیگران

انواع کپچا را به شکل‌های مختلفی می‌توان دسته‌بندی کرد. در یک دسته‌بندی کپچاها به دو دسته کپچاهای مبتنی بر OCR و کپچاهای غیر مبتنی بر OCR دسته‌بندی می‌شوند. در دسته‌بندی دیگر کپچاها به ۳ دسته تقسیم می‌شوند که در ادامه به بررسی هر یک از آن‌ها می‌پردازیم و کارهای انجام شده در هر بخش را معرفی می‌کنیم.

۲-۱- کپچاهای مبتنی بر کاراکتر

در این نوع کپچا رشته‌ای از کاراکترها به کاربر نمایش داده می‌شود که این رشته می‌تواند شامل لغات یا ترکیبی تصادفی از کاراکترها باشد [۱۰]. در زمینه کپچاهای مبتنی بر کاراکتر کارهای فراوانی صورت گرفته است که از آن جمله می‌توان به EZ-Gimpy دانشگاه کارنگی ملون که توسط شرکت Yahoo استفاده می‌شد، اشاره کرد [۱۱].

در این کپچا از روش اعوجاج برای تفکیک کاربران انسانی از نرم‌افزارهای رایانه‌ای استفاده می‌شد اما چون در این روش کلمات از یک فرهنگ لغت با ۸۵۰ کلمه انتخاب می‌شدند، به راحتی در برابر حملات شکسته شد [۱۱].



شکل ۱: Yahoo's EZ-Gimpy [۱۲]

نمونه‌ای دیگر از کارهای انجام شده کپچای Gimpy است که در آن کاربر باید ۳ تا از ۷ کلمه نمایش داده شده در تصویر را تشخیص دهد [۱۱]. نمونه‌ای از این کپچا در شکل ۲ آمده است.



شکل ۲: نمونه کپچای Gimpy

۳-۱- ویژگی‌های الفبای فارسی [۱]

- شکل‌های مختلف حروف: الفبای فارسی شامل ۳۲ حرف است که هرکدام از آن‌ها با توجه به مکانشان در یک کلمه می‌توانند به ۲ تا ۴ شکل مختلف ظاهر شوند.
- اندازه متفاوت حروف فارسی نسبت به یکدیگر متصل بودن حروف
- نقطه‌دار بودن حروف: بسیاری از حروف فارسی مشابه هستند و تنها تفاوتشان در یک یا چند نقطه و یا علائمی خاص است.
- امکان کشیدگی بعضی از حروف نوشتن از راست به چپ
- وجود علائم خاص مثل تشدید، تنوین و مد: به کاربرد این علائم باعث مشکل کردن کار OCR در تشخیص نویز از این علائم می‌گردد.

- تغییر ضمیر متصل متناسب با فاعل: در زبان فارسی ضمیر متصل به فعل با توجه به فاعل جمله تغییر می‌کند و یا در برخی حالات حذف می‌شود. مثلا مصدر "رفتن" در صورتی که با اول شخص مفرد به کار رود به صورت "رفتم" و در صورتی که با دوم شخص مفرد به کار رود به صورت "رفتی" صرف می‌شود که این ویژگی از ویژگی‌های زبان فارسی است و در زبان لاتین به این صورت وجود ندارد. از این ویژگی می‌توان برای ایجاد کپچای کارآمدتر با قابلیت اطمینان بیشتر استفاده کرد.

۳-۲- ساختار جملات زبان فارسی

با توجه به اینکه در طرح ارائه شده در این مقاله از جملات فارسی برای ساخت کپچا استفاده شده است. در این بخش به بررسی ساختار جملات فارسی می‌پردازیم.

به طور کلی جملات زبان فارسی براساس تعداد اجزای اصلیشان به سه دسته تقسیم می‌شوند:

۱. جملات دو جزئی: جملاتی هستند که فعل آن‌ها ناگذر است. یعنی فقط به نهاد نیاز دارد. ساختار جملات دوجزئی:

نهاد + فعل

۲. جملات سه جزئی: جملاتی هستند که علاوه بر نهاد به مفعول، متمم یا مسند هم نیاز دارند. ساختار انواع جملات سه جزئی در زیر آمده است:

(۱) نهاد + مفعول + فعل، (۲) نهاد + متمم + فعل

(۳) نهاد + مسند + فعل

۳. جملات چهار جزئی: این جملات با چهار ساختار زیر ساخته می‌شوند:

Please select all the cat photos:



شکل ۵: نمونه کپچای Asirra شرکت Microsoft

Bongo کپچای دیگری است که دو مجموعه از اشکال هندسی را به کاربر نمایش می‌دهد. پس از اینکه کاربر این دو مجموعه را دید یک شکل هندسی به عنوان سوال به کاربر نشان داده می‌شود که کاربر باید تشخیص دهد این شکل مربوط به مجموعه اشکال هندسی سمت چپ است یا در مجموعه سمت راست وجود دارد.

ESP-Pix دانشگاه کارنگی ملون نمونه دیگری از کپچاهای مبتنی بر تصویر است که چند تصویر با موضوع مشابه را به کاربر نمایش می‌دهد و از کاربر می‌خواهد که مشخص کند، تصاویر در چه زمینه‌ای هستند [۵]. نمونه ای از این نوع در شکل ۶ آمده است.



شکل ۶: [۱۲] CMU's ESP-PIX

از دیگر کارهای انجام شده می‌توان به کپچای IMAGINATION [۱۲] و همچنین کپچای تصویری دانشگاه برکلی [۱۸] اشاره کرد.

۳-۲- کپچاهای مبتنی بر صوت

در این نوع کپچا کاربر به یک فایل صوتی گوش می‌دهد و کاراکترها یا اعدادی را که شنیده است به عنوان پاسخ کپچا وارد می‌کند. برای جلوگیری از شکسته شدن این کپچا توسط نرم‌افزارهای تشخیص گفتار، از تکنیک اعوجاج بر روی صوت استفاده می‌شود که این کار تشخیص صوت را برای انسان نیز مشکل می‌کند [۲۰]. reCAPTCHA که در بخش ۲-۱ به آن اشاره شد، نمونه‌ای از این نوع کپچا است که در دسته‌بندی کپچاهای مبتنی بر کاراکتر نیز قرار می‌گیرد. در زمینه کپچاهای مبتنی بر صوت کارهای متنوع دیگری نیز انجام شده است [۲۱ و ۲۰].

۳- ویژگی‌های الفبا و ساختار جملات فارسی، کپچای فارسی

در کپچای ارائه شده از ویژگی‌های ساختاری الفبای فارسی برای بالا بردن قابلیت اطمینان کپچا استفاده شده است. از این رو در این بخش ابتدا به طور مختصر به ویژگی‌های زبان فارسی اشاره می‌کنیم و پس از آن ساختار جملات فارسی را بررسی می‌کنیم و در نهایت به توضیح کارهای انجام شده در زمینه کپچای فارسی می‌پردازیم.

(۱) نهاد + مفعول + متمم + فعل، (۲) نهاد + مفعول + مسند + فعل

(۳) نهاد + متمم + مسند + فعل، (۴) نهاد + مفعول + مفعول + فعل

۳-۳- کپچای فارسی

در این بخش به معرفی کارهای انجام شده در زمینه کپچای فارسی می‌پردازیم. از جمله اولین کارهای انجام شده در این زمینه کپچای Persian/Arabic BaffleText توسط شهرضا [۷] می‌باشد.

در این کپچا یک کلمه بی معنی فارسی یا عربی که بین ۳ تا ۸ کاراکتر دارد به صورت تصادفی ساخته شده و بر روی یک پس‌زمینه رنگی قرار می‌گیرد و برای بالا بردن کارایی کپچا خطوطی به صورت تصادفی به پس‌زمینه تصویر اضافه می‌شوند (شکل ۷). در واقع این کپچا از برخی ویژگی‌های زبان فارسی از جمله نقطه‌دار بودن حروف و متصل بودن حروف به یکدیگر استفاده کرده‌است همچنین در ساخت کلمات این کپچا سعی شده‌است از حروف مشابه فارسی مثل "ب، پ، ت، ث" در کنار هم استفاده شود تا تشخیص کلمه برای OCR سخت‌تر شود. از دیگر ویژگی‌های این نمونه استفاده از قلم تصادفی برای ساخت کلمه نمایش داده شده به کاربر است.



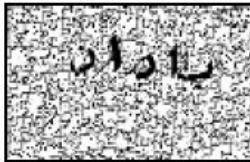
شکل ۷: نمونه کپچای شهرضا [۷]

کپچای نستعلیق فارسی [۲۲] نمونه دیگری از کپچاهای فارسی است. در این کپچا که باز هم توسط شهرضا طراحی شده از خط نستعلیق و کلمات با معنی فارسی استفاده شده‌است. تصاویر از یک بانک اطلاعاتی استخراج شده و به کاربر نمایش داده می‌شوند و هیچ‌گونه اعوجاجی به تصویر داده نشده‌است که این کار باعث کاهش کارایی این کپچا می‌شود، زیرا اخیراً روش‌هایی جهت تشخیص کلمات نستعلیق ارائه شده‌است [۲۳]. نمونه‌ای از تصاویر به کار رفته در این کپچا در شکل ۸ آمده‌است.



شکل ۸: نمونه کپچای نستعلیق فارسی [۲۲]

از دیگر نمونه‌های کپچای فارسی کپچای طراحی شده توسط بهلول و ملک‌زاده [۲] می‌باشد. در این کپچا یک کلمه بی‌معنی فارسی ساخته شده و بر روی تصویر قرار می‌گیرد و برای جلوگیری از تشخیص آن توسط OCR تغییراتی از جمله جابه‌جایی، خمیدگی، شیف‌ت و خط‌خطی کردن بر روی تصویر اعمال می‌شود. در نهایت برای افزایش ضریب اطمینان کپچا نویز نمک-فلفلی نیز مطابق شکل ۹ به تصویر اضافه می‌شود.



شکل ۹: نمونه کپچای بهلول و ملک‌زاده [۲]

نمونه‌ای دیگر از کپچای فارسی در [۱] آمده‌است. این کپچا از کلمات معنادار فارسی و نمایش آن به صورت کلمات بدون نقطه در تصویر بهره می‌گیرد. به این صورت که سه کلمه بدون نقطه در یک تصویر با پس‌زمینه رنگی قرار می‌گیرد و سپس بر روی تصویر نویز اعمال می‌شود و تصویر نهایی به کاربر نمایش داده می‌شود. برای اعمال نویز در این کپچا از یک سری خطوط و بیضی استفاده شده‌است. نمونه‌ای از این کپچا در شکل ۱۰ مشاهده می‌شود.



شکل ۱۰: نمونه کپچای [۱]

از دیگر کارهای انجام شده در این زمینه می‌توان به کپچای تصویری مرجع [۳] اشاره کرد که در شکل ۱۱ آمده‌است.



شکل ۱۱: نمونه کپچای [۳]

۴- الگوریتم پیشنهادی

استفاده از جملات فارسی به صورتی که فاعل یا فعل جمله نمایش داده شده از لحاظ دستور زبان فارسی صحیح نباشد روشی است که برای ایجاد نوعی از کپچای فارسی در این مقاله مورد استفاده قرار گرفته‌است. به این صورت که یک جمله به کاربر نمایش داده می‌شود و از او خواسته می‌شود که فعل جمله یا فاعل آن را اصلاح کند.

به عنوان مثال در شکل‌های ۱۲ و ۱۳ نمونه‌هایی از این کپچا برای درک بهتر آن آمده‌است.



شکل ۱۲: نمونه کپچای ارائه شده (اصلاح فاعل)



شکل ۱۳: نمونه کپچای ارائه شده (اصلاح فعل)

در شکل ۱۲ از کاربر خواسته شده تا فاعل جمله را اصلاح کند. همانطور که مشاهده می‌نمایید جمله "شما دیروز به کارخانه رفتیم" از لحاظ فاعل (شما) و فعل (رفتیم) مطابقت ندارد که در اینجا کاربر باید کلمه "من" را به عنوان پاسخ وارد کند. در شکل ۱۳ جمله "تو شش ماه قبل پول را می‌پردازی" به کاربر نمایش داده شده است ولی این بار از کاربر خواسته شده تا فعل جمله را اصلاح کند. با توجه به قید زمان جمله (شش ماه قبل) که یک قید زمان گذشته است، انتظار می‌رود فعل جمله نیز گذشته باشد ولی فعل جمله بالا (می‌پردازی) فعل مضارع است و با زمان جمله مطابقت ندارد. در اینجا هریک از پاسخ‌های "پرداختی"، "پرداخته‌ای" و یا "پرداخته‌بودی" که همگی بر زمان گذشته دلالت دارند، قابل قبول است.

انتظار می‌رود که به دلیل تاکید بر اصلاح جمله از نظر دستوری طرح ارائه شده از قابلیت اطمینان نسبتاً بالا برخوردار باشد. چرا که ماشین درک صحیحی از معنای کلمات ندارد و حتی در صورتی که OCR بتواند جمله ساخته شده در تصویر را تشخیص دهد، امکان عبور از کپچا برای ماشین وجود نخواهد داشت. چرا که ماشین باید بتواند جمله را از نظر دستوری اصلاح کند و برای این کار نیاز به درک قیود زمان، فاعل و فعل جمله خواهد داشت که با توجه به گستردگی آن‌ها به راحتی امکان پذیر نیست. از سوی دیگر اصلاح فعل و فاعل جملات فارسی برای کاربران فارسی زبان به سهولت امکان پذیر است.

۴-۱-۱-۱-۴ - گردآوری کلمات مناسب

برای عملی ساختن طرح ارائه شده، سیستم باید بتواند جملات بامعنای فارسی را بسازد. برای اینکه از پیچیده شدن جملات کپچا جلوگیری شود و کار اصلاح جملات برای کاربران به سادگی امکان پذیر باشد، در این کپچا از جملات دو جزئی، سه جزئی با مفعول و سه جزئی با متمم استفاده و از جملات سه جزئی با مسند و چهار جزئی به طور کلی صرف نظر شده است. برای ساخت جملات نیاز به جمع‌آوری کلمات فارسی و دسته‌بندی آن‌ها براساس نقششان در جمله می‌باشد. از این رو پایگاه داده‌ای از کلمات فارسی با توجه به نقش کلمه در جمله ساخته شد. در این پایگاه داده کلمات به چهار دسته "فعل"، "مفعول-متمم"، "فاعل" و "قید" تقسیم‌بندی شده‌اند.

۴-۱-۱-۲-۱-۴ - فعل

در بخش فعل تعدادی از افعال متداول فارسی که در جملات دوجزئی، سه جزئی با مفعول و یا سه جزئی با متمم به کار می‌روند، جمع‌آوری شد و افعال براساس نوع جملاتی که در آن‌ها به کار می‌روند، دسته‌بندی شدند. لازم به ذکر است که افعال به صورت بن‌ماضی همراه با بن‌مضارع متناسب در پایگاه داده ذخیره شدند تا بتوان جملاتی با زمان‌های مختلف ایجاد کرد. علاوه بر بن‌ماضی و مضارع، برای فعل‌های گذرا به متمم، حرف اضافه متناسب با فعل نیز در پایگاه داده ذخیره شد.

۴-۱-۲-۲-۱-۴ - مفعول-متمم

در این بخش کلماتی از دسته‌بندی‌های متفاوت شامل اسامی حیوانات، اشیاء، وسایل منزل، انواع لباس، غذاها و همچنین کلمات کاربردی دیگر در زمینه‌های اجتماعی، اقتصادی، موسیقی و ... جمع‌آوری و در پایگاه داده ذخیره شد. کلمات این بخش طوری انتخاب شدند که با افعال ذخیره شده در پایگاه داده متناسب باشند.

۴-۱-۲-۳-۱-۴ - فاعل

در بخش فاعل ۶ ضمیر "من، تو، او، ما، شما و آن‌ها" به همراه تعداد زیادی از نام‌های فارسی در پایگاه داده ذخیره شدند.

۴-۱-۲-۴-۱-۴ - قید

تعدادی از قیود زمان که بر زمان گذشته، حال یا آینده دلالت می‌کنند به همراه زمان آن‌ها در پایگاه داده ذخیره شد. مثال‌هایی از قیود ذخیره شده در پایگاه داده: امسال، ماه آینده، دیروز، اکنون، فردا

۴-۱-۲-۴-۲-۴ - ساخت جمله

پس از جمع‌آوری کلمات در چهار بخش ذکر شده در بخش قبل، سیستم باید بتواند جملات صحیح فارسی از نظر دستوری را با استفاده از پایگاه داده بسازد. برای این کار باید مراحل زیر طی شوند:

۱. انتخاب یک فاعل به صورت تصادفی و مشخص کردن شخص آن
۲. انتخاب تصادفی یک قید زمان و مشخص کردن زمان یا زمان‌هایی که آن قید در آن‌ها می‌تواند مورد استفاده قرار گیرد.
۳. انتخاب یکی از زمان‌های مجاز به صورت تصادفی
۴. انتخاب تصادفی یک فعل از میان افعال ذخیره‌شده و تغییر آن با توجه به شخص فاعل و زمان انتخاب شده در مرحله ۳
۵. بررسی فعل از لحاظ نوع جملاتی که در آن‌ها کاربرد دارد.
۶. انتخاب تصادفی مفعول یا متمم، متناسب با نوع تشخیص داده‌شده در مرحله ۵
۷. ساخت جمله بر اساس ساختار جملات زبان فارسی

۴-۱-۲-۳-۲-۴ - به هم ریختن جمله از لحاظ دستوری

در این بخش جمله ساخته شده در بخش قبلی به گونه‌ای تغییر می‌کند که جمله حاصل، از لحاظ دستور زبان فارسی با مشکل مواجه شود. سیستم به صورت تصادفی یکی از دو بخش "زمان و شخص فعل" و یا "فاعل" جمله را با توجه به "زمان و شخص فعل"، "زمان‌های مجاز برای قید زمان" و "شخص فاعل" به هم می‌ریزد. سپس از کاربر خواسته می‌شود که همان بخش را اصلاح کند تا جمله از لحاظ دستوری درست شود.

۴-۱-۲-۴-۲-۴ - ساختار و محل قرارگیری جمله در تصویر

سیستم پس از به هم ریختن جمله براساس طول جمله اندازه قلم را طوری تعیین می‌کند که جمله به طور کامل در تصویر قرار گیرد. به منظور

نتایج بدست آمده از بین ۴۳۰ کاربر که به کپچا پاسخ داده‌اند با توجه به اینکه در بار اول، دوم، سوم یا دفعات بعدی به کپچا پاسخ صحیح داده‌اند در شکل ۱۵ آمده‌است.



شکل ۱۵: درصد صحیح پاسخ‌گویی به کپچا

با توجه به نتایج بدست آمده مجموعاً ۸۵ درصد کاربران در بار اول یا دوم به کپچا پاسخ صحیح داده‌اند که نشان دهنده این است که اگرچه کپچای طراحی شده در برابر OCR از امنیت بالایی برخوردار است، حل آن برای انسان خیلی مشکل نیست.

با توجه به آمار ارائه شده و مقایسه با سایر کارهای انجام شده در زمینه کپچای فارسی، پروژه حاضر جزء کارهای قدرتمند در زمینه کپچا محسوب می‌شود. ضمن اینکه مجموعاً ۷۸ درصد کاربران از کپچای ارائه شده اظهار رضایت کرده‌اند.

۶- نتیجه‌گیری

در این مقاله به معرفی کارهای انجام شده در زمینه کپچا پرداختیم. با وجود کارهای فراوان و متنوع انجام شده در کپچای انگلیسی، هنوز در کپچای فارسی کار زیادی انجام نشده‌است و با وجود ویژگی‌های منحصر به فرد زبان فارسی که به برخی از آن‌ها اشاره شد کپچای فارسی جای کار بسیار زیادی دارد.

در مقاله حاضر سعی شد تا با استفاده از ساختار دستور زبان فارسی کپچای فارسی جدید، امن و باکارایی بالا طراحی شود. استفاده از جملات زبان فارسی و به هم ریختن فعل و فاعل جمله از مواردی است که برای افزایش امنیت کپچا در نظر گرفته شده‌است. به طوری که یک کاربر فارسی زبان به راحتی معنای جمله را درک کرده و آن را اصلاح می‌نماید ولی نرم‌افزارهای تشخیص متن به راحتی قادر به اصلاح جمله نیستند.

در پایان می‌توان گفت که به نظر می‌رسد پروژه حاضر توانسته باشد هدف اولیه خود یعنی ایجاد کپچای فارسی قدرتمند با قابلیت اطمینان و کارایی بالا را تحقق بخشیده باشد. ادامه کار می‌تواند شامل مواردی چون جابه‌جا کردن اجزای جمله و یا اضافه کردن یک جزء جدید به جملات (مثلاً قید مکان) باشد. به نحوی که تشخیص جمله برای ماشین مشکل‌تر شود ولی کاربران قادر به تشخیص جمله درست باشند.

افزایش سطح امنیت کپچا زاویه و محل قرار گیری جمله از نظر افقی و عمودی در تصویر ثابت نیست بلکه به صورت تصادفی تغییر می‌کند به گونه‌ای که در تمامی حالات، کل جمله در تصویر قرار گیرد و از آن خارج نشود. در این صورت تشخیص اجزای جمله برای OCR مشکل‌تر می‌شود.

۴-۵- تولید تصویر اولیه

در این مرحله تنظیمات و چیدمانی که در مرحله قبل برای نمایش جمله در نظر گرفته شده‌است، با استفاده از توابع گرافیکی به تصویر اولیه تبدیل می‌شود. طول و عرض تصویر ایجاد شده ۱۰۰*۴۰۰ پیکسل است که می‌توان با تغییر اندازه قلم آن را کوچکتر یا بزرگتر کرد.

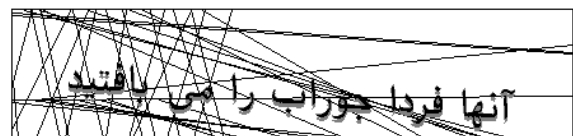
۴-۵-۱- افزودن سایه به تصویر

برای اینکه کار ماشین برای تشخیص اجزای جمله کمی مشکل‌تر شود، سایه‌ای با فاصله کم و به رنگ خاکستری به جمله ایجاد شده در تصویر اضافه می‌نماییم.

۴-۵-۲- افزودن نویز به تصویر

در مورد زبان فارسی، اضافه کردن نویز، شباهت‌هایی را با علائم و نقطه‌های کلمات فارسی به وجود می‌آورد، که کار را برای نرم‌افزارهای تشخیص متن مشکل می‌کند [۱].

برای افزودن نویز، یک سری خطوط به تصویر ساخته شده در مرحله قبل اضافه می‌نماییم. تعداد، زاویه و محل قرارگیری خطوط به صورت تصادفی و به گونه‌ای انتخاب می‌شوند که تصویر برای برنامه‌های کامپیوتری مبهم و برای انسان واضح باشد. نمونه‌ای از تصویر نهایی در شکل ۱۴ مشاهده می‌شود.



شکل ۱۴: نمونه تصویر نهایی (آنها فردا جوراب را می‌بافتند)

۵- نتایج پیاده‌سازی

روش ارائه شده در این مقاله به زبان PHP پیاده‌سازی شد و جهت آزمایش توسط کاربران بر روی یک سایت در محیط وب قرار گرفت و از افراد بالای ۱۵ سال خواسته شد تا کپچا را حل نمایند. به منظور ارزیابی نتایج، پایگاه داده‌ای برای ذخیره‌سازی اطلاعات ایجاد گردید که در آن IP کاربر، جمله نمایش داده شده به کاربر، پاسخی که به آن داده‌است، تاریخ و ساعت وارد کردن کپچا ذخیره و از اطلاعات فوق برای سنجش امنیت کپچای طراحی شده استفاده شد. همچنین از کاربران خواسته شد تا نظر خود را در مورد این کپچا با انتخاب یکی از گزینه‌های "عالی"، "خوب"، "متوسط" و "بد" بیان کنند.

Mellon University, Science, Vol 321, no. 5895, pp. 1465-1468, 2008.

[15] "reCAPTCHA: Stop Spam Read Books", Retrieved on 24 December 2013, at <http://www.recaptcha.net>

[16] A. L. Coates, H. S. Baird, and R. J. Fateman, "Pessimistic Print: a reverse Turing Test", in Proc. of the 6th International Conf. on Document Analysis and Recognition, Seattle, US, pp. 1154-1158, Sep. 2001.

[17] M. Chew, and H. S. Baird, "BaffleText: a human interactive proof", in Proc. of the 10th SPIE/IS&T Conf. on Document Recognition and Retrieval (DR&R2003), Santa Clara, US, pp. 305-316, Jan. 2003.

[18] M. Chew, and J. D. Tygar, "Image Recognition CAPTCHAs", in Proc. of the 7th International Information Security Conference (ISC 2004), Springer, pp. 268-279, September 2004.

[19] J. Elson, J. Douceur, and J. Saul, "Asirra: A CAPTCHA that exploits Interest-Aligned Manual Image Categorization", Proc. of the 14th ACM conference on Computer and communications security, 2007.

[20] A. Schlaikjer, "A Dual-Use Speech CAPTCHA: Aiding Visually Impaired Web Users While Providing Transcriptions and Audio Streams", Language Technologies Institute School of Computer Science, Carnegie Mellon University, CMU-LTI-07-014., <http://www.lti.cs.cmu.edu>, November 2007.

[21] J. Holman, J. Lazar, J. Feng, and J. D'Arcy, "Developing Usable CAPTCHAs for Blind Users", Proc. of the 9th international ACM SIGACCESS conference on Computers and accessibility, 2007.

[22] M. H. Shirali-Shahreza, M. Shirali-Shahreza, "Nastaliq CAPTCHA", Iranian Journal of Electrical and Computer Engineering (IJECE), Vol. 5, No. 2, pp. 109-114, 2007.

[23] M. Salmani Jelodar, M. J. Fadaeieslam, N. Mozayani, M. Fazeli, "A Persian OCR System Using Morphological Operators", Transactions on Engineering, Computing and Technology, V4 February 2005 ISSN 1305-5313, Manuscript received January 21, 2005.

زیر نویس ها

¹ CAPTCHA: Completely Automated Public Turing test to tell Computers and Humans Apart

² OCR: Optical Character Recognition

³ Animal Species Image Recognition for Restricting Access

[۱] یغمایی، فرزین؛ کامیار، محدثه؛ کمندی، فائزه، "طراحی و پیاده سازی کپچای فارسی قدرتمند با نمایشی بدون نقطه از کلمات"، بیست و یکمین کنفرانس مهندسی برق ایران، دانشگاه فردوسی مشهد، اردیبهشت ۱۳۹۲

[۲] بهلول، مهدی؛ ملک‌زاده، محمد، "سیستم کپچای فارسی برای جلوگیری از ثبت نام خودکار ربات‌های نرم افزاری در صفحات وب"، سیزدهمین کنفرانس انجمن کامپیوتر، جزیره کیش، سال ۱۳۸۶

[۳] یغمایی، فرزین؛ بخشنده، عبدالجبار، "روشی جدید برای تولید کپچای فارسی با مقاومت بالا مقابل حملات"، نهمین کنفرانس بین المللی انجمن رمز ایران، شهریور ۱۳۹۱

[4] I. Zeifman, "The Incapsula Blog Report : Bot traffic is up to 61.5 % of all website traffic", Retrieved on 24 December 2013, at <http://www.incapsula.com/the-incapsula-blog>

[5] L. Von Ahn, M. Blum, and J. Langford, Telling Human and Computers Apart Automatically, Communications of the ACM, February 2004, 47(2), pp. 57-60.

[6] M. Blum, L. Von Ahn, and J. Langford, "Completely Automatic Public Turing Test to Tell Computers and Humans Apart", the CAPTCHA Project, Department of Computer Science, Carnegie-Mellon University, 2000, Retrieved on 15 December 2013, at <http://www.captcha.net>

[7] M. H. Shirali-Shahreza, M. Shirali-Shahreza, "Persian/Arabic BaffleText CAPTCHA", Journal of Universal Computer Science, vol. 12, no. 12(2006), 1783-1796.

[8] E. Walsh, "CAPTCHA cracked by artificial intelligence", 28 October 2013, Retrieved on 27 November 2013, at <http://mybroadband.co.za/news/internet/90435-captcha-cracked-by-artificial-intelligence.html>

[9] J. Newman, "Ads Within Captchas: Tell Me it Isn't So", 23 September 2010, Retrieved on 16 December 2013, at http://www.pcworld.com/article/206051/Ads_Within_Captchas_Tell_Me_it_Isnt_So.html

[10] G. Sauer, H. Hochheiser, J. Feng, and J. Lazar, "Towards a Universally Usable CAPTCHA", Department of Computer and Information Sciences, Towson University, In Proc. of the 4th Symp. On Usable Privacy and Security (SOUPS), 2008.

[11] G. Mori, J. Malik, "Recognizing Objects in Adversarial Clutter: Breaking a Visual CAPTCHA", Computer Science Division University of California, Berkeley, Proc. of IEEE CS Society Conf. on Computer Vision and Pattern Recognition (CVPR'03), Madison, WI, 2003, pp. 131-141

[12] R. Datta, J. Li, and J. Z. Wang, "Exploiting the Human-Machine Gap in Image Recognition for Designing CAPTCHAs", IEEE Transactions on Information Forensics and Security, Vol 4, Issue:3, pp. 504-518, 2009

[13] H. S. Baird, T. Riopka, "ScatterType: a Reading CAPTCHA Resistant to Segmentation Attack", Computer Science & Engineering, Lehigh University, Accepted for publication in Proceedings, IS & T/SPIE Document Recognition & Retrieval XII Conference, San Jose, CA, January 16-20, 2005.

[14] L. Von Ahn, B. Maurer, C. McMillen, D. Abraham, M. Blum, "reCAPTCHA: Human-Based Character Recognition via Web Security Measures", Computer Science Department, Carnegie