

OCR: آموزش الفبای فارسی به رایانه

رایانه تنها یک ابزار است، و ما انسان‌ها می‌باید شیوه کار کردن را به او بیاموزیم، و امکانات کار را هم برایش فراهم آوریم، نرم‌افزارهای مختلف به همین منظور طراحی می‌شوند. یکی از آنها نرم‌افزار OCR است که برای تشخیص و بازیابی الفبا (نوشته‌های دست‌نویس یا تایپ‌شده) توسط کامپیوتر طراحی می‌شوند.

اگر هر یک از ما به طور متوسط ۷۰ تا ۸۰ سال عمر کنیم مجموع روزهای زندگیمان بین ۲۵ هزار تا ۲۹ هزار روز می‌شود. حال اگر به ما بگویند نرم‌افزاری تولید می‌شود که می‌تواند بسیار بیشتر از عمر چندین و چند انسان در وقت صرفه‌جویی کند حتماً از آن استقبال می‌کنیم، مگر آنکه برای چند صد برابر روزهای عمرمان ارزشی قائل نباشیم!

نرم‌افزار OCR قادر است چنین تحویلی را در استفاده وقت بشر به وجود آورد. این نرم‌افزار می‌تواند متن‌های دست‌نویس یا متونی را که قبلاً تایپ شده‌اند و اکنون فایل تایپی آنها موجود نیست، خود به تنهایی و بدون دخالت انگشتان هیچ تایپیستی تایپ کند. فرض کنید که مثلاً می‌خواهیم متن مقالات روزنامه اطلاعات سال ۱۳۴۰ شمسی را (که اکنون نه تنها فایل تایپی‌اش موجود نیست - چون آن زمان اصلاً تایپ کامپیوتری در کار نبود! - بلکه خود نسخه‌های روزنامه را هم به زحمت می‌توان پیدا کرد) تایپ دیجیتالی کنیم، و این متن‌ها را داخل بسته‌های نرم‌افزاری یا اینترنت قرار دهیم. اگر هر شماره از روزنامه را ۲۴ صفحه فرض کنیم، و هر تایپیست بتواند در هر روز حداکثر یک صفحه از آن صفحات کاهی و کهنه شده قدیمی را دوباره تایپ کند، مجموعاً ۲۴ روز لازم است تا تنها مقالات یک شماره از روزنامه تایپ شود. بنابراین در عرض یک سال یک نفر می‌تواند تنها ۱۵ شماره از روزنامه را تایپ کند. حال اگر نرم‌افزاری باشیم که بتواند با اسکن کردن هر صفحه روزنامه، به طور خودکار مقالات آن را تایپ کند، تحویلی عظیم رخ می‌دهد، یعنی مطالب و مقالات هزاران شماره از روزنامه‌های قدیمی به سرعت وارد فایل‌های رایانه‌ای می‌شود. حال این امکان را تعمیم بدهید به هزاران کتاب و دست‌نویس‌های قدیمی یا جدید، که هر کس بخواهد تنها یک صفحه از آنها را تایپ کند، باید کلی وقت صرف کند. می‌بینید که نرم‌افزار OCR برآستی می‌تواند هزاران هزار روز در وقت ما صرفه‌جویی کند، و البته هزینه‌ها را هم کاهش دهد. البته فقط یک مشکل کوچک به وجود می‌آید و آن بیکار شدن تایپیست‌هاست!

OCR در ایران چگونه آغاز شد؟

امان از دست این تیزهوشان! ماجرا از ثبت‌نام داوطلبان آزمون «سازمان ملی پرورش استعدادهاي درخشان (تیزهوشان)» در سال ۱۳۸۰ آغاز شد. ثبت‌نام از روی فرم‌هایی که توسط دانش‌آموزان تکمیل می‌شد انجام می‌گرفت. دانش‌آموزان شرکت‌کننده در آزمون - مانند آزمون‌های سراسری - باید نام، نام خانوادگی، نام پدر، نام شهرستان محل تولد و سکونت، نام مدرسه و دین خود را در داخل کادرهای مربعی شکل و به صورت حروف مقطع (یعنی هر حرف داخل یک کادر) می‌نوشتند. وقتی که همه فرم‌ها از طریق پست به سازمان مرکزی برگزارکننده آزمون می‌رسید، عده زیادی تایپیست متن آنها را دوباره وارد رایانه می‌کردند. در واقع همان حرف‌های داخل کادر را دوباره تایپ می‌کردند تا اطلاعات شناسنامه‌ای هر دانش‌آموز به صورت دیجیتالی درآید. این روش هم بسیار زمان‌بر بود و هم نیاز به تعداد زیادی تایپیست داشت. احتمال داشت که تایپیست‌ها هم هنگام تایپ اشتباه کنند و با ثبت نادرست یک نام، مشخصات فردی در رایانه مرکزی وارد شود که اصلاً متولد نشده است! مثلاً فرض کنید تایپیست محترم نام «جواد» را، که داخل کادرها به صورت «ج.و.ا.د» نوشته شده بود، «فؤاد» تایپ می‌کرد؛ در آن صورت در کارت شناسایی

جواد سابق، فؤاد فعلی ثبت می‌شد! (جواد موجود حذف می‌شد و فؤاد ناموجود وارد فهرست داوطلبان می‌شد!) افزون بر این، هزینه کار نیز بسیار زیاد بود. به علت همین مشکلات، در بهمن‌ماه ۱۳۸۰، نخستین طرح OCR برای بازشناسی حروف فارسی توسط کامپیوتر ارائه شد و در سال‌های ۱۳۸۱ و ۱۳۸۲ نیز ثبت‌نام آزمون تیزهوشان به یاری این نرم‌افزار انجام شد.

OCR چیست؟

OCR سرنام اصطلاحی است که صورت کامل آن در واژه‌نامه‌های انگلیسی به دو صورت آمده است:

۱. Recognition Optical Character

۲. Optical Character Reader

فرض کنید که ما متنی را روی کاغذ داریم و می‌خواهیم آن را وارد رایانه کنیم. اولین روشی که به ذهن می‌رسد این است که متن را به تایپست بدهیم تا با کامپیوتر تایپ کند. اما آیا می‌شود عین همان متن را وارد رایانه بکنیم تا نیازی به تایپ نباشد؟ البته دستگاه «اسکنر» می‌تواند تصویری از آن متن را وارد رایانه کند، تا اینجا بخشی از مشکل ما حل شده است. اما رایانه که نه عقلی دارد و نه «زبان» می‌فهمد، نمی‌تواند حروف و کلمات را از هم تشخیص دهد. مثلاً اگر از کامپیوتر بخواهیم به ما بگوید که در متن اسکن‌شده کلمه «علی» چند بار آمده است، بی‌آنکه شرمنده شود، می‌گوید: «error»، یعنی: «نمی‌توانم تشخیص بدهم!» در واقع این «تصویر دیجیتالی‌شده» باید به «تصویر قابل پردازش» تبدیل شود. موضوع اصلی OCR همین است.

انواع OCR

در زبان‌های دیگر، به ویژه زبان‌هایی که با حروف لاتینی نوشته می‌شوند، سال‌هاست که از OCR استفاده می‌شود. اما در ایران تازه دو سه سالی است که به فکر استفاده از OCR در زبان فارسی افتاده‌ایم.

و اما OCR چند نوع است: یا تایپی است یا دست‌نویس. یعنی یا باید یک متن قبلاً تایپ شده را (مثل کتاب‌ها و روزنامه‌های چندین سال قبل، یا حتی متنی را که فایل تایپی آن موجود نیست و فقط پرینت آن را داریم) وارد رایانه کنیم، یا متن دست‌نویس را. متن‌های دست‌نویس هم به دو صورت «گسسته» و «پیوسته» وجود دارند: متن «دست‌نویس پیوسته» مثل همان چیزهایی است که ما هرازگاهی که دلمان تنگ می‌شود روی کاغذ می‌نویسیم، یا یک نامه، یا یک قطعه شعر و ... اما متن «دست‌نویس گسسته» همان نوشته‌هایی است که حروف آن جدا از هم و به صورت گسسته نوشته شده‌اند، مثل نام و نام‌خانوادگی که در فرم‌های آزمون ثبت‌نام، به صورت هر حرف داخل یک کادر، نوشته می‌شوند. طراحی OCR گسسته فارسی تقریباً در مراحل پایانی کار قرار دارد ولی، OCR پیوسته ظاهراً سال‌های زیادی کار می‌برد.

«رضا صدیق» و «پرویز رزازی»، که در رشته مخابرات تحصیل کرده‌اند و مسئولان یک شرکت کامپیوتری به نام «اندیشه نرم‌افزار پایا» هستند، برای اولین بار به طور جدی پروژه OCR فارسی را دنبال کرده‌اند. رزازی که دانشجوی مخابرات و مسئول بخش پردازش سیگنال شرکت «پایا» و مدیر پروژه OCR در این شرکت است، می‌گوید: «OCR در دنیا موضوعی ناشناخته نیست، و بر روی آن زیاد کار شده است، ولی در ایران با آنکه مدت‌هاست روی آن کار شده، اما بسیاری از این کارها در حد کارهای دانشگاهی و مقاله‌های علمی باقی‌مانده بود و تبدیل به یک محصول کاربردی در ابعاد وسیع (مثل ثبت‌نام آزمون‌های بزرگ) نشده بود. ما بر روی این طرح کار کردیم و هدفمان هم این بود که محصول را به شکل صنعتی آن تولید کنیم.»

البته غیر از شرکت «پایا»، دو شرکت دیگر نیز با حمایت دبیرخانه طرح «تکفا» (توسعه کاربرد فناوری اطلاعات و ارتباطات) مشغول پژوهش و آزمایش بر روی OCR فارسی هستند. یکی از این شرکت‌ها «داده‌پردازان دوران نوین» نام دارد که مدیریت آن را دکتر «حسام فیلی» بر عهده دارد. دکتر فیلی متخصص در رشته هوش مصنوعی، از دانشگاه صنعتی شریف، است و شرکت «دوران نوین» را از سال ۱۳۸۱، با هدف کار تخصصی بر روی پروژه‌های هوش مصنوعی تأسیس کرده است. او درباره چگونگی پیوستن شرکتش به این طرح می‌گوید: «از تیرماه سال ۸۲ با شروع فعالیت طرح «تکفا» و حمایت‌های مالی آنها، این شرکت تصمیم گرفت که در زمینه طراحی OCR فارسی پژوهش و فعالیت کند. این پروژه در شرکت «دوران نوین» با همکاری آقای «دکتر ابراهیمی مقدم» که او هم از دانشجویان دوره دکتری هوش مصنوعی دانشگاه صنعتی شریف است، انجام می‌گیرد.

فارسی ما و مشکلات آن

قبل از اینکه به مراحل دیگر OCR پردازیم، لازم است اندکی هم به مشکلات خط فارسی – یا در واقع ویژگی‌های این خط – پردازیم. اول اینکه ما در فارسی حروف را به صورت چسبیده و پیوسته می‌نویسیم و این کار برای تشخیص حرف به حرف نوشته از سوی رایانه (که قرار است در مراحل بعدی آن را تایپ کند) بسیار مشکل است. تصور کنید که همین کلمه ساده «است» را به حالت‌های مختلف می‌شود نوشت: یکی برای «س» دندان می‌گذارد، یکی نمی‌گذارد، یکی آن را می‌کشد و یکی نمی‌کشد و... حالا اگر همین صورت‌های مختلف «س» به «ت» هم بچسبند، تشخیص حروف برای ما انسان‌ها هم سخت می‌شود، چه رسد به رایانه.

شباهت حروف

مشکل دیگر خط ما این است که حرف‌های فارسی بسیار به هم شبیه‌اند. مثلاً در نظر بگیرید که تفاوت «ر» با «ز» یا «ذ» یا «ب» با «ت» تنها در يك نقطه است، و چون نقطه جزء بسیار کوچکی است، اگر يك خط یا حتي يك لك كوچك روي كاغذ بیفتد، تشخیص حروف از هم بسیار دشوار می‌شود و دردسر جدي برای بازشناسی حروف توسط رایانه ایجاد می‌کند. اینها تازه مشکلات خط فارسی است. درباره اعداد فارسی هم این مشکل وجود دارد: صفر ما يك نقطه كوچك است که می‌تواند رایانه را به اشتباه بیندازد؛ اعداد ۱، ۲، ۳ هم بسیار به هم شبیه هستند و تنها تفاوتشان يك دندانه كوچك است.

به دلایل گفته شده OCR در مرحله کنونی در کشور ما مربوط به «دست‌نویس‌های گسسته» یا متن‌های تایپی پیوسته است، و تا بازشناسی متن‌های دست‌نویس پیوسته توسط کامپیوتر راه زیادی در پیش است، چون در دست‌نویس‌های گسسته، اگرچه حروف به هم شباهت دارند، حداقل جداجدا نوشته شده‌اند. در متن‌های پیوسته تایپی هم مشکل کشیده شدن يك حرف یا شکسته نوشته شدن حروف را نداریم. البته به گفته مسئولان شرکت «پایا» در حال حاضر هم نرم‌افزارهایی وجود دارد که متن دست‌نویس پیوسته را تبدیل به حروف جدا از هم و گسسته می‌کنند، ولی ضریب خطای این نرم‌افزارها زیاد است و به شکل صنعتی درنیامده‌اند.

بازشناسی حروف و الگو

تا اینجا گفتیم تصویر صفحه‌ای که در آن حروف به طور جداجدا (هر حرف داخل يك کادر) نوشته شده است، به وسیله اسکن وارد رایانه می‌شود. مرحله بعدی این است که حروف بازشناسی شوند، یعنی مکان آنها از دیگر خطوط (مثل خطوط کادری که داخل آن نوشته شده) بازشناسی شود، و اگر متن پیوسته تایپی است، حروف جدا شوند و زواید تصویر حذف شود. مثلاً اگر دانش‌آموزی «س» را به گونه‌ای نوشت که بیرون از کادر بود، به رایانه بفهمانیم که بی‌دقتی شده است او باید همان حرف داخل کادر را بخواند. در

مرحله بعدی که «بازشناسی الگو» نام دارد، با تعدادی شرط می‌شود فهمید که مثلاً حرفی «الف» است یا نه، و رایانه تشخیص می‌دهد که حرف «پ» است یا «ب». برای این تشخیص لازم است که تصویر حرف «الف» با الفهای نمونه - که قبلاً به رایانه داده شده است - منطبق شود. الفبای نمونه قبلاً از روی یک مجموعه بزرگ آموزشی تهیه شده و ویژگی‌های مشترک از آن استخراج شده است. اما از آنجا که تنوع صورت‌ها نوشتاری یک حرف به صورت دست‌نویس بسیار زیاد است، مدلی آماری استخراج می‌شود که در آن شباهت ویژگی‌های استخراج شده قبلی با نمونه ورودی به رایانه بررسی می‌شود. در اینجا «بازشناسی الگو» با روش‌های آماری انجام می‌شود که روش معمول در سیستم‌های OCR است.

اگر فکر می‌کنید که کار تمام شده است در اشتباهید، چون تازه می‌رسیم به دنباله حروف. مثلاً اگر کسی همان حرف «س» را با دنباله بنویسد، رایانه باید تشخیص دهد که این حرف فقط «س» است، یا مثلاً «ی» هم به آن چسبیده است.

مدلسازی یا پردازش زبانی

مرحله بعدی «مدلسازی زبانی» یا «پردازش زبانی» نام دارد. حروف به هم چسبیده، که کلمه را درست می‌کنند، باید معنی‌دار یا شناخته‌شده باشند. در این مرحله بررسی می‌شود که چه کلماتی در زبان وجود دارد؟ چه ترکیب‌هایی از کلمات مجاز است؟ و... البته در مراحل پیشرفته‌تر، مدلسازی گرامری (دستور زبان) و مدلسازی معنایی هم وجود دارد که تشخیص می‌دهد جمله از لحاظ دستوری و معنایی درست است یا بی‌مفهوم است. اما در OCR گسسته - که بیشتر برای ثبت نام استفاده شده - شباهت یک کلمه به نام، نام خانوادگی، شهر و ... کافی است.

برای تشخیص ترکیب‌های مجاز یک کلمه یا معنی‌دار بودن یک کلمه نیز به تهیه بانک‌های اطلاعاتی (Data base) نیاز داریم. در این بانک‌ها مثلاً تمام نام‌های کوچک و بزرگ ایرانیان قبلاً جمع‌آوری شده است و هنگام تطبیق یک کلمه با آن مشخص می‌شود که رایانه حروف آن را دست تشخیص داده یا نه. بنابراین نقش این بانک اطلاعاتی بسیار مهم است، چون اگر نامی در آن ثبت نشده باشد، کلمه‌ای که آن نام را شامل شود، به طور خودکار از برنامه OCR حذف می‌شود یا پیغام می‌آید که: «این کلمه اشتباه است» در صورتی که ممکن است مثلاً نام «هشام» در بین نام‌های ایرانی وجود داشته باشد، ولی قبلاً در بانک اطلاعاتی ثبت نشده باشد.

بانک‌های ما و دیگران

مهندس «رزازی» درباره مشکل بانک‌های اطلاعاتی در زبان فارسی می‌گوید: «در دنیا برای توسعه OCR و ارزیابی آن، بانک‌های اطلاعاتی استاندارد ساخته شده است که در آنها همه کلمات وجود دارند، یعنی بانک هم مشکل دیجیتال کلمه را دارد، و هم تصویرش را. اما برای زبان فارسی، این بانک‌های اطلاعاتی چه برای ارزیابی و چه برای توسعه، استاندارد شده نیست. در واقع هر کسی برای خودش یک بانک اطلاعاتی می‌سازد، و این نمونه‌های متفاوت مشکلاتی را ایجاد می‌کند. مثلاً برای ثبت نام دانش‌آموزانی که در آزمون مدارس تیزهوشان شرکت کرده بودند، یک بانک اطلاعاتی حاوی نام‌های فارسی، از روی اطلاعات فرم‌های سال‌های قبل، تهیه شد که از روی آن کلماتی که خیلی شبیه به نام‌های فارسی بودند تشخیص داده می‌شد. مثلاً اگر رایانه کلمه‌ای را «مصیبت» تشخیص داد، براساس بانک اطلاعاتی معلوم می‌شود که «مصیب» بوده است که یک نام ایرانی است.

علي، ولي، قلي ... و سيب نکته ديگر اين است كه يك بانك اطلاعاتي بايد شامل تعدادي كلمات خام باشد، بلكه «بسامد» آن واژگان، يعني ميزان استعمال و تكرار كلمات در زبان و مشخصات آمري آنها هم بايد ثبت شده باشد، والا كارايي زيادي ندارد. مثلاً «علي» نامي است كه شباهت زيادي به «ولي» و «قلي» دارد. كارهاي آمري در بانك اطلاعاتي بايد طوري انجام شده باشد كه تعداد «علي» بيشتر باشد، تا و بعد نوبت «ولي» و «قلي» برسد، چون درصد بسامدي «علي» به لحاظ آمري و کاربرد در ميان نامهاي بيشتر است. در OCR فارسي گسسته، اگر فقط مربوط به نامها و نامخانوادگي باشد، كار ساده تر است از حالي كه در OCR پيوسته وجود دارد. چون در OCR پيوسته هر كلمه اي ممكن است وجود داشته باشد مثل «سيب»، اما در بانك اطلاعاتي نامها همه مي دانيم كه سيب نام يك شخص نيست بلكه نام يك ميوه است! بنا بر اين در OCR همواره سعي مي شود كه درصد خطا کاهش يابد، تا كلمات در حد ممكن درست تشخيص داده شوند. اگرچه طراحان هنوز به صد درصد صحت نرسيده اند، ولي نگران نتايج آزمون خود نباشيد، چون تمامي اطلاعات مربوط به شما چندين بار كنترل مي شوند و از سازوكار reject (يا مردودي) در رايانه هم استفاده مي شود. در اين روش اگر رايانه نتوانست كلمه اي را تشخيص دهد، مي فهمد كه نفهميده است و در خروجي اش مي آورد كه: «من اين كلمه را نفهميده ام» و كار به سيستم دستي مي رود و در آنجا تصحيح مي شود. اين فرآيند در پست خيلي كارايي دارد. در هر جاي دنيا كه تفكيك نامه ها و ديگر مرسولات پستي به وسيله OCR انجام مي شود، بعضي از نامه ها در سيستم كامپيوتري وارد سازوكار «مردودي» مي شوند و به طور دستي مورد بررسي مجدد قرار مي گيرند. هم اكنون در سطح محدودتي از OCR در پست كشور ما نيز استفاده مي شود، چون در پست هم كد پستي چند رقمي و ديگر اطلاعات به صورت گسسته و داخل كادرهاي نوشته مي شود، و كار آسان تر است.

در مورد خطاي OCR در تشخيص كلمات، مسئولان شركت «پايا» نظر جالبي دارند:

«حتي با تعبيه سيستم مردودي (reject) هم ممكن است خطايي در تشخيص كلمات وجود داشته باشد. بايد در نظر داشته باشيم كه هيچ سيستم پردازشگري (از جمله انسان) بدون خطا نيست. نکته مهم اين است كه يك سيستم ماشيني درصد خطاي كم تري نسبت به انسان داشته باشد تا جايگزين خوبي براي انسان باشد. مسئله اين نيست كه خطا را به صفر برسانيم. هر قدر كه فناوري جلوتر مي رود، ميزان خطا هم بيشتر کاهش مي يابد.»

مدیر شرکت «دوران نوین» هم به گونه اي ديگر به همين موضوع اشاره مي کند: «انتظار ما از مسئولان طرح «تكفا» آن است كه با موضوع OCR واقع بينانه تر برخورد شود، و در بحث مربوط به هزينه هاي پروژه و انتظاراتي كه از OCR مي رود، واقعيتها در نظر گرفته شود. ديده كنوني مسئولان تكفا آن است كه كل مشكل «خطا» تا ۱۰۰ درصد حل شود، در حالي كه فكر مي كنم حل مسائل مربوط به هوش مصنوعي نياز به روش تدريجي دارد. مثلاً در زبان عربي هم، نرم افزار «صخر» در نسخه اول خود فقط تا حدود ۴۰ درصد دقت داشت، در حالي كه اكنون پس از گذشت ۱۲ سال از اولين نسخه آن دقت به مرز ۹۸ درصد رسيده است.»

همان طور كه اشاره شد از OCR در ثبت نام آزمون «سازمان ملي استعدادهاي درخشان» در سالهاي ۸۱ و ۸۲ استفاده شد كه از طريق آن ۴۴۰،۰۰۰ نفر به طور ماشيني ثبت نام شدند. اين روش باعث شد كه در سال ۸۱ (نمونه اول) ۴۵ درصد در هزينه ها و ۲۵ درصد در زمان ثبت نام صرفه جويي شود. در سال بعد (۸۲) اين رقم به ۵۰ درصد رسيد.

نرم‌افزاری که در این آزمون‌ها مورد استفاده قرار گرفت برای هر کدام از موارد صحت بازشناسی متفاوتی داشت و در مجموع کار آن خوب بود. به نظر می‌آید که در چند سال آینده و با پیشرفت OCR فارسی و کاهش هر چه بیشتر خطای آن، در آزمون‌های بزرگ‌تری مانند آزمون سراسری دانشگاه‌ها نیز بتوان از آن استفاده کرد.

سرنوشت OCR دست‌نویس

در مورد OCR پیوسته دست‌نویس نیز روند کار به همان صورتی است که شرح دادیم، اما آنچه کار را دشوارتر می‌کند، قطعه‌بندی و جداجدا کردن حروف به هم چسبیده و تشخیص آنهاست. اگر این روند طی شود، این امید وجود دارد که روزی از OCR پیوسته دست‌نویس فارسی هم در سطح گسترده‌ای استفاده شود. البته OCR پیوسته دست‌نویس حتی در زبان انگلیسی هم هنوز به کاربرد وسیع و عملی نرسیده است. مهندس «رزازی» در این مورد می‌گوید: «OCR انگلیسی در سیستم عامل windows وجود دارد که همراه با office فروخته می‌شود، ولی فکر نکنید که نامه‌های اداری انگلیسی که با دست‌نویس نوشته شده‌اند همه با OCR تایپ می‌شوند. این کار برای به نتیجه رسیدن به حداقل یک روند ۱۰ ساله را باید طی کند. OCR فارسی یک مرحله عقب‌تر است، پس زمان بیشتری می‌برد.» مهندس «صدیق»، مدیرعامل شرکت «پایا» هم می‌گوید: «همین OCR فارسی گسسته هم تا چند سال پیش یک رؤیا بود، ولی دیدیم که محقق شده است و به مرور پیشرفته‌تر هم خواهد شد. بنابراین طراحی OCR پیوسته فارسی هم، اگرچه سال‌ها طول می‌کشد، ولی مطمئناً به نتیجه خواهد رسید. این طرح یک طرح تحقیقاتی است که در دانشگاه‌ها دنبال می‌شود و هنوز به یک محصول صنعتی قابل استفاده در سطح کلان و کاربردی برای عموم نرسیده است. ولی در حال حاضر نمونه‌های دانشگاهی و آزمایشگاهی آن در داخل کشور وجود دارد و موضوع رساله دکتری برخی از دانشجویان است.»

بنابراین بین ۱۰ تا ۲۰ سال آینده، آن‌گونه که مسئولان شرکت «پایا» می‌گویند، OCR پیوسته دست‌نویس فارسی هم وارد بازار خواهد شد. دکتر فیلی هم در پاسخ به این سؤال که «آیا طراحی OCR پیوسته فارسی روزی تحقق خواهد یافت؟» پاسخ می‌دهد: «بله ولی به تدریج.» به هر حال براساس قرارداد «تکفا» با شرکت‌های ایرانی، تا کمتر از یک ماه دیگر، نسخه‌نهایی (البته نه صددرصد تکمیل‌شده) OCR فارسی دست‌نویس گسسته و تایپی پیوسته ارائه خواهد شد. مدیر شرکت «دوران نوین» در این مورد می‌گوید: «پروژه OCR گسسته در مراحل پایانی خود قرار دارد ولی دارای مشکلاتی در تشخیص انواع اسکنرها و انواع فونت‌هاست که در حال رفع آن هستیم. این نرم‌افزار در حال حاضر امکان تشخیص فونت‌های تایپی فارسی با دقت زیاد را دارد، ولی مشکل جدی آن است که با اسکنرهای مختلف نتایج نامناسبی می‌دهد.» وی از اهمیت این طرح در بعد کلان ملی هم می‌گوید: «با توجه به این که مشکل OCR برای بسیاری از زبان‌های دنیا مانند انگلیسی عملاً حل شده است، اگر در کشور ما هم به نتیجه‌نهایی برسد در افزایش سطح اطلاعات فارسی در دنیای دیجیتال امروز (از جمله در اینترنت) بسیار اهمیت خواهد داشت.»



mohammad6347@yahoo.com