



فرامرز هنروران^۱

او.سی. آر و کاربردهای آن در کتابخانه‌ها و مراکز اطلاع‌رسانی



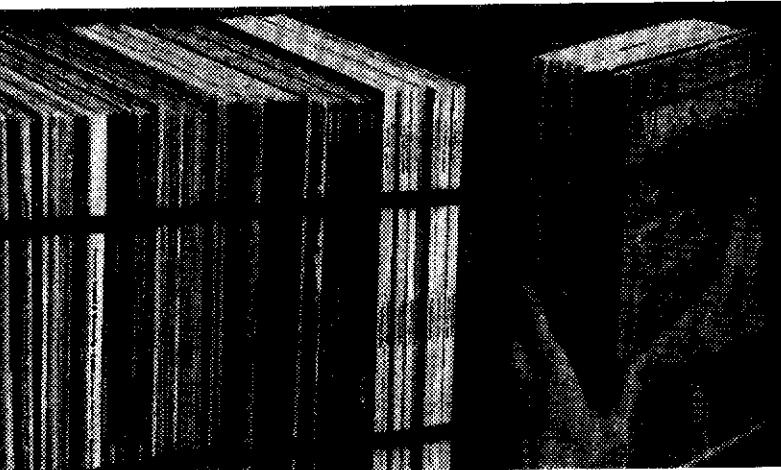
متخصصان ساخت افزاری و نرم افزاری کامپیوتر و دستگاه‌های جانبی آن به فکر تهیه دستگاه‌ها و برنامه‌های باشند تا نیاز به صفحه کلید کامپیوتر را برای ورود اطلاعاتی که قبل از روی کاغذ نقش بسته بر طرف سازد. به غیر از عامل بر شمرده شده دو دلیل دیگر وجود و توسعه چنین برنامه‌هایی را تقویت می‌کند: ۱) حجم اطلاعاتی که بر روی کاغذ از زمان اختراع چاپ ناکنون درج شده نسبت به اطلاعاتی که بر روی حافظه‌های ثانویه کامپیوتر (دیسک و دیسکت، نسوار مغناطیسی...) موجود است از میزان بسیار بالاتر برخوردار است و اگر قصد ورود تمامی این اطلاعات به کمک کاربران باشد انجامش به سالها وقت و تعداد بیشتر کامپیوتر نیاز داشته و هزینه فراوانی می‌طلبد؛ ۲) اشتباهاتی که توسط کاربر در هنگام ورود اطلاعات رخ می‌دهد در بیشتر موارد به ۸ درصد یا بالاتر می‌رسد و زمانی که جهت رفع این اغلاط تلف می‌شود بسیار زیاد است. وجود مسائل بالا منجر به تولید و تکمیل دستگاه‌های گوناگون پویشگر^۲ و نرم افزارهای مختلف شناسائی بصری و هوشمند حروف گردیده است. پویشگر وسیله‌ای است که یک صفحه چاپی - چه متنی و چه تصویری، یا ترکیبی از هر دو - را پویش می‌کند. طریقه عمل بدین صورت است که ابتدا نوری به صفحه چاپی از طریق یک منبع نوری (که بر حسب نوع پویشگر به صورت خط به خط یا صفحه به صفحه است) تابانیده شده و نور منعکس شده از صفحه به طرف آرایه‌ای از دیودهای نوری^۳ هدایت می‌شود. این نور پس از برخورد با این دیودها نوسانات الکتریکی تولید می‌کند که متنطبق با نقاط سیاه (یا رنگی) و زمینه کاغذ که تشکیل دهنده متن و یا تصویرند، می‌باشد این نوسانات به صورت قیاسی^۴ به کمک مدارات الکترونیکی پویشگر تقویت شده و تبدیل به اطلاعات رقمه‌ای^۵ که به صورت آرایه‌ای از اطلاعات دودویی و یا نقشه بیتی در آمده و به حافظه RAM^۶ (حافظه با دسترسی دلخواه) یا حافظه ثانویه کامپیوتر منتقل می‌گردد و در این مرحله است که کار نرم افزار او سی آر شروع می‌شود. این نرم افزارها نقشه بیتی تولید شده توسط پویشگر را خوانده و در حین فرآیندی خاص قسمتهای سفید یا خالی صفحه پویش شده را مشخص می‌کند تا بدین طریق قادر شود محل خطوط متن و تصاویر احتمالی موجود را برای ادامه کار

چند سالی از بحث کتابخانه‌های بدون کتاب و اداره‌های بدون کاغذ^۷ می‌گذرد ولی با نگاهی اجمالی به وضعیت فعلی در دنیا به سادگی می‌توان متفاوت شد که هنوز راهی طولانی برای رسیدن به این مهم به طور نسبتاً آرمانی باقی مانده است. کاغذ به سبب سهولت استفاده و سندبیش هنوز بکی از محمل‌های اصلی اطلاعات به شمار می‌رود. ولی حجم زیاد، عمر محدود، کندی بازیابی و انتقال آنچه بر این محمل نقش بسته استفاده از آن را در عصری که دستیابی و انتقال هر چه سریعتر اطلاعات را عامل تعیین کننده کرده است، چه در زمینه درج اطلاعات جدید و چه تکثیر اطلاعات قدیمی روزی روز غیرااقتصادی تر می‌سازد، از جهت دیگر تولید آن را با توجه به محدود بودن منابع گیاهی و وارد آوردن خسارات جبران ناپذیر به محیط زیست، هر روز به نقطه‌ای بحرانی تر نزدیک می‌سازد. مجموع این عوامل باعث شده که طی بیست و چند سال گذشته

OCR

در مرحله ماقبل آخر که در بعضی نرم افزارهای او سی آر وجود دارد با کمک یک برنامه غلطیاب املایی، کلمات استخراج شده، بررسی شده و اغلاط احتمالی در موقع شناسایی کاراکترهای هر کلمه بدینوسیله تصحیح می شود. در پایان نرم افزار او سی آر تمامی کاراکترهای مدرک را به صورت فایل متنی اسکن^{۱۶} یعنی به شکلی که در کامپیوتر قابل استفاده است تبدیل کرده و در این حال است که می توان متن را به نرم افزارهای واژه پرداز یا نرم افزارهای بانک اطلاعاتی متنی داد یا آن را به هر محل دیگر از طریق کامپیوتر منتقل کرد، تمام

خود تعیین کند. در مرحله اول تبدیل تصویر به متن، این نرم افزار سعی می کند تا با مقابله نقطه به نقطه^۹ تصویر صفحه پوشش شده سند با مجموعه کاراکترهای تعریف شده که نرم افزار او سی آر به حافظه RAM کامپیوتر سپرده است کاراکترهای کامل^{۱۰} مشابه را شناسایی کرده و از تصویر استخراج نماید. این مجموعه کاراکترهای تعریف شده شامل تمامی فونتهای پرکاربرد (fonnt^{۱۱}) مجموعه ای است از حروف چاہی در اندازه، شکل و سبک مشخص تشکیل شده از حروف بزرگ و کوچک، اعداد و علامت مختلف) الفبای لاتین (يونانی یا سریلیک) است. بدین سبب که در این مرحله از شناسایی حروف انطباق کامل لازم است اگر جگونگی تصویر پوشش شده از لحاظ کیفیت پایین باشد این روش چندان کارآمد نبوده و باید از روش دیگری بهره برد که در مرحله بعدی موجب شناسایی کاراکترهای شناخته شده در مرحله اول می شود. این روش وقتیگر به نام استخراج ویژگی^{۱۲} یا تحلیل ویژگی^{۱۳} یکی از روشهایی است که علاوه بر نرم افزارهای او سی آر در نرم افزارهای ICR^{۱۴} (نرم افزارهایی برای استخراج حروف از متن خطی یا چاپ شده با فونتهای غیراستاندارد) نیز کاربرد وسیعی دارد. در این روش ابعاد، نسبت ابعاد حواشی، تعداد گوشه ها و لبه های هر کاراکتر مورد تجزیه و تحلیل قرار گرفته و نشایه آن با کاراکترهای تعریف شده سنجیده می شود، به عبارت دیگر، در این روش ارتفاع هر کاراکتر نسبت به خط کرسی محاسبه شده و هر کاراکتر به صورت ترکیبی از خطوط مستقیم و منحنی های باز و بسته در مقایسه با یکایک کاراکترهای تعریف شده مورد تجزیه و تحلیل قرار می گیرد تا از نتایج این تجزیه و تحلیل کاراکتر شناسایی شود. در مرحله سوم به دلیل آن که دو مرحله قبل امکان شناسایی تمامی کاراکترهای متن را نمی دهد، نرم افزار به دو روش (بسته به نوع نرم افزار) برای کشف کاراکترهای ناشناخته کمک می طلبد، یا علامتی همچون @ یا # را جایگزین کاراکترهای ناشناخته می کند و به کاربر اعلام می دارد که خود این علامت را پیدا کرده و کاراکتر مناسب را جایگزین این علامت سازد و یا: هر کاراکتر ناشناخته را بر روی صفحه نمایش کامپیوتر^{۱۵} مشخص کرده و از کاربر می خواهد تا کلید مربوط به این کاراکتر را بر روی صفحه کلید کامپیوتر فشار دهد تا به جای کاراکتر مجهول کاراکتر صحیح جانشین شود.



مراحل بالا بسته به نوع مدرک و نوع نرم افزار و ساخت افزار به کار رفته از یک تا چند دقیقه وقت می برد که در مقایسه با انجام این کار با کمک کاربر در اکثر موارد از سرعتی ۱۰ تا ۳۰ درصد بالاتر برخوردار است، و ذکر این نکته نیز خالی از فایده نیست که بر حسب تجربیاتی که انجام گرفته مشخص شده که هزینه پیدا کردن و تصحیح یک حرف (البته در الفبای لاتین) حدود ده هزار برابر هزینه تصحیح یک حرف رد شده بدلیل ناخوانای بودن یا غیرقابل تشخیص بودن از طرف نرم افزار او سی آر است و همچنین تجربیات مشخص ساخته که استفاده از او سی آر حدود ۵ درصد در هزینه ها صرفه جویی می کند.

او سی آر در کتابخانه ها و مراکز اطلاع رسانی کتابخانه ها و مراکز اطلاع رسانی از جمله مؤسسانی هستند که تهیه اطلاعات و انتقال آن اصلی ترین عامل وجودی شان بوده و هر چه در این امور سریعتر عمل کنند به هدف شان نزدیکتر می شوند. در کتابخانه های گشتویی بویزه



به کار مشغول است دست به انتخاب یکی از این نرم‌افزارها زده یا کلاً^{۱۸} با مزایا و معایب هر یک از آنها آشنا شود. این نرم افزارها عبارتند از:

Caere Omnipage Professional 5.0.۱
Recognita Plus 2.0 International ۲. از شرکت Recognita

Xerox Imaging Systems TextBridge 2.0.۳
Calera Recognition WordScan Plus ۴. از شرکت Systems

این مقایسه‌ها از چندین جنبه توسط هوارد اگلواشتاين^{۱۸} صورت گرفته و در مجله بایت مورخ اکتبر ۱۹۹۴ درج گردیده که پرداختن به تمامی آنها از حوصله این مقاله خارج است و تنها به ذکر نکات اساسی این مقایسه پرداخته می‌شود. دو عامل اساسی در مقایسه نرم‌افزارهای او سی آر صحبت انجام کار و زمان نسبی انجام کار این نرم‌افزارهای است. صحبت انجام کار او سی آر عبارت است از حاصل تقسیم تعداد حروفی که به طور صحیح شناسایی شده بر تعداد کل حروف در یک سند بر حسب درصد. نتایج حاصل از مقایسه صحبت کار این نرم‌افزارها در جدول ۱ مشخص است، که حاصل او سی آر ۳۰ صفحه آ است.

زمان نسبی انجام کار نرم‌افزارهای او سی آر یا به بیان دیگر سرعت انجام کار، عبارت است از حاصل تقسیم تعداد کلی حروفی که به طور صحیح او سی آر شده به زمانی که صرف شناسایی حروف کل سند می‌شود. کلاً در نرم‌افزارهای او سی آر صحبت از سرعت مهمن تر است. در این مقایسه نرم‌افزار TextBridge از سرعت بالا و صحبت معقولی برخوردار بوده و WordScan و OmniPage از صحبت بالا برخوردار بوده و Recognita سریع بوده ولی خطای زیاد داشته است. مزیت‌های دیگری که یک نرم‌افزار او سی آر می‌تواند داشته باشد عبارتند از:

۱. راحتی کارکردن با نرم‌افزار؛
۲. قابلیت کار با دستگاه‌های مختلف پوششگر؛
۳. تشخیص فونتهای مختلف؛
۴. قابلیت شناسایی حروف زبانهای مختلف؛
۵. قیمت مناسب.

که از میان این نرم‌افزارها، TextBridge راحت‌ترین کاربرد

کتابخانه‌های کشورهای در حال توسعه همه چیز با لاقل درصد اطلاعات موجود بر کاغذ نقش بسته است. از فهرست‌برگه تا یکایک نامه‌های بخش‌های مختلف یک کتابخانه - همان گونه که گفته شد - از مهمترین قابنیت‌یعنی امکان دستیابی و انتقال سریع بی‌بهراه‌اند علاوه بر آن امکان جستجو، تصحیح، تجزیه و تحلیل آماری در آنها تنها به روش دستی امکان‌پذیر است. اما با استفاده از روش او سی آر تمامی این اطلاعات در زمانی کوتاه و با هزینه‌ای پایین نسبت به روش‌های دیگر خواهد توانست از تمامی قابنیت‌های ذکر شده بهره‌مند گردد. به طور مثال برای ارائه متن مقالات مجلات و روزنامه‌ها، دیگر نیازی به جستجو در میان صدھا برگه در برگه‌دانها یا جستجو در میان صدھا تصویر میکروfilm، نخواهد بود. زیرا با استخراج متن این مقالات توسط او سی آر و دادن این اطلاعات به برنامه‌های بانک اطلاعات متنی، به راحتی می‌توان به هر مقاله دست یافت و آن را به کمک خطوط مخابراتی به هر نقطه‌ای فرستاد؛ و صدھا مورد استفاده دیگر که بسته به نوع کتابخانه و میزان وسعت آن متفاوت است.

معرفی چهار برنامه مشهور او سی آر و مقایسه آنها با هم

حال که در مورد مزایای روش او سی آر و کاربردهای آن در حوزه کتابداری و اطلاع‌رسانی بحث شد، به معرفی و مقایسه چهار برنامه او سی آر که تحت برنامه ویندوز مایکروسافت^{۱۷} کار می‌کند، پرداخته می‌شود تا بتوان در ضمن این مقایسه با توجه به نیازهای کتابخانه یا مؤسسه‌ای که در آن

جدول ۱. مقایسه صحبت کار او سی آر در چهار نرم‌افزار*

	سن نهیه شده توسط چاپگر چرخ خورشیدی	Omnipage	Recognita	TextBridge	Wordscan
سن نهیه شده توسط چاپگر چرخ خورشیدی	٪۹۹/۳	٪۹۷/۰	٪۹۹/۲	٪۹۸/۹	٪۹۸/۹
سن نهیه شده توسط چاپگر چوهرافشان لیزری	٪۹۹/۳	٪۹۶/۲	٪۹۸/۷	٪۹۹/۰	٪۹۸/۰
من کمی شده ^{**}	٪۸۸/۹	٪۹۲/۳	٪۹۶/۴	٪۹۵/۲	٪۹۸/۰
من فاکس شده ^{**}	٪۹۸/۸	٪۸۷/۴	٪۷۸/۱	٪۹۸/۰	٪۹۸/۰

* عدم صحبت ۱ درصدیدان معنی است که در هر ۱۰۰ حرف بک حرف اشتباه شناسایی شده است.

** من اصلی این دو مورد به توسط چاپگر لیزری با وضوح بالا نهیه شده است.

مجلات کامپیوتربازی داخلی گاه و بیگانه خبرهایی هر چند کوتاه و نیمه موشک در انجام این تلاش به اطلاع علاقمندان و متخصصان می‌رسد که امید است این تلاشها به شر برسد و بزوی شاهد نزم افزایی مناسب برای خط فارسی باشیم.

و کمترین قیمت را داشته و Recognita از بیشترین قابلیت در شناسایی حروف زبانهای اروپایی یا نوشته شده با خط لاتین و یونانی برخوردار بوده (حدود ۸۰ زبان) و در ضمن قابلیت کار باشیش از ۱۰۰ نوع دستگاه پویشگر را دارد (این نتایج باتوجه به داده‌های جدول ۱ حاصل شده است).

١٧

1. Schantz, Herbert F. "Using OCR: Three Crucial Variables". in: *Inform*, 6 (June 1992) PP. 20-22.
 2. Kahama, Paz Y. "Forms Removal: Saving Storage Space, improving OCR Performance". ibd. PP 12-18.
 3. Longley, Denis & Shain, Michael, *Dictionary of Information technology*. London: Macmillan, 1989.
 4. *Encyclopedia of Computer Science and Engineering*. New York: Van Nostrand Reinhold Co. 1983.
 5. Eglowstein, Howard, "Due Recognition for OCR" in: *Byte*, October 1994. PP. 145-148.

یادداشت‌ها:

۱. کارشناس ارشد کتابداری و اطلاع‌رسانی دانشکده روانشناسی و علوم تربیتی دانشگاه تهران.

- 1- Optical Character Recognition (OCR)
 - 2- Paperless Office
 - 3- Scanner
 - 4- Photodiodes' Matrix
 - 5- Analog
 - 6- Digital
 - 7- Bitmap
 - 8- Random Access Memory
 - 9- Pixel by Pixel
 - 10- Font
 - 11- Feature Extraction
 - 12- Feature analysis
 - 13- Intelligent Character Recognition
 - 14- Monitor
 - 15- Keyboard
 - 16- ASCII (American Standard Code for Information Interchange)
 - 17- MS Windows
 - 18- Howard Feltwein



18