

SISTEM PENTRU ASISTAREA INTRODUCERII ÎN CALCULATOR A DATELOR DIN CÂMPURILE FORMULARELOR TIPIZATE

Mihnea Horia Vrejoiu

mihnea@dossvl.ici.ro

Institutul Național de Cercetare-Dezvoltare în Informatică - ICI București

Rezumat: Introducerea manuală a datelor de pe suport de hârtie în calculator reprezintă o activitate pe cât de necesară și importantă în contextul informatizării pe scară largă în mai toate domeniile, pe atât de consumatoare de resurse umane și de timp și o potențială sursă de erori de tastare, mai ales în cazul volumelor mari de date ce trebuie preluate. În acest context, sunt extrem de utile tehnici și instrumente prin care această activitate poate fi automatizată. Adesea însă, datorită cerințelor critice privind acuratețea de preluare a datelor este neapărat necesară măcar o verificare și validare a datelor de către operatorul uman în lipsa altor criterii și/sau posibilității de validare automată. În acest context s-a propus și implementat experimental o soluție software semiautomată, bazată pe tehnici OCR/ICR, destinată asistării activității de introducere în calculator a datelor completate în formulare tipizate, cu format fix, pentru creșterea eficienței, productivității și acurateței. Au fost efectuate experimente de testare, au fost făcute observații asupra funcționării și rezultatelor obținute și au fost sintetizate câteva concluzii și posibilități de îmbunătățiri și optimizări ulterioare.

Cuvinte cheie: introducere de date în calculator, formulare tipizate, învățare automată, învățare supervizată, OCR/ICR, expresii regulate.

Abstract: Manually inputting data from paper support represents, on one hand, a necessary and important activity in the context of the widely spread informatization in most areas, but, on the other hand, so human, and time, resources consuming, and also a potential source of typing errors, especially when large volumes of data are inputted. In this context, techniques and tools by which this activity can be automated are extremely useful. However, often, due to critical requirements on the accuracy of the inputted data, it is absolutely necessary at least a check and validation of these data by the human operator in the absence of other criteria and/or automated validation possibilities. In this context, it has been proposed and experimentally implemented a semi-automatic software solution based on OCR/ICR techniques, intended to assist the work of inputting the data filled in standardized forms, with fixed format, for improving the efficiency, productivity and accuracy. Testing experiments were conducted, there were made observations on the functioning and results, there were summarized some conclusions and possible further improvements and optimizations.

Keywords: data inputting, standardized forms, machine learning, supervised learning, OCR/ICR, regular expressions.

1. Introducere

În pofida evoluției tehnologice și a introducerii informatizării pe scară largă în cele mai multe sectoare de activitate, precum și a politicilor și eforturilor ecologiste de protejare a mediului și de sustenabilitate actuale, totuși volumul de informație vehiculată pe suport de hârtie și care trebuie introdusă în calculator este încă foarte mare.

Un exemplu este cel al formularelor tipizate. Preluarea manuală în calculator a datelor din acestea poate implica volume mari de timp și efort și poate introduce erori de tastare, în unele situații inacceptabile. Automatizarea procesului a reprezentat una din problematicile cele mai vizate de marii producători de software în domeniul analizei de imagini OCR/ICR și management al documentelor. Există deja soluții comerciale pe piață pentru preluarea "automată" a datelor din formulare, cu funcțiuni și performanțe diferite, dar și altele anunțate care nu s-au impus (încă) pentru o utilizare pe scară largă. În continuare prezentăm succint câteva idei generale, sintetice, privind acest tip de aplicații.

Cele mai multe soluții sunt configurabile și se bazează pe crearea și folosirea de șabloane / machete asociate tipurilor de formulare care specifică poziționarea și caracteristicile câmpurilor acestora. Unele își propun chiar o identificare automată a câmpurilor de interes pe diferite criterii de analiză a imaginii formularului scanat.

Majoritatea produselor de acest tip anunță capacitatea de recunoaștere atât a caracterelor tipărite cât și a celor scrise de mână (separate între ele). De asemenea, unele produse permit în plus și recunoașterea marcajelor din check-box-uri și a codurilor de bare;

Deși majoritatea producătorilor anunță rate de recunoaștere spectaculoase ale motoarelor

OCR/ICR, toate aplicațiile acordă o mare atenție etapei de corecție și validare post-recunoaștere. Corecțiile post-recunoaștere se fac în general prin utilizarea de dicționare, informații apriori despre câmpuri, calcularea unor sume de control etc., în timp ce validarea necesită întotdeauna intervenția utilizatorului.

Astfel, se poate afirma că soluțiile existente sunt practic semiautomate, ele oferind instrumente și interfețe de asistare a operatorului uman în această activitate.

Majoritatea produselor oferă facilități de conectare, în ceea ce privește exportul de date, cu diverse sisteme de baze de date standardizate și/sau proprietare.

2. Sistemul software propus

S-a propus și implementat experimental o modalitate semiautomată de asistare a operatorului uman pentru creșterea eficienței, productivității și acurateții în activitatea de introducere în calculator a datelor completate din câmpurile formularelor tipizate, cu format fix. Aceasta se realizează prin identificarea și preluarea automată utilizând tehnici de recunoaștere optică / inteligentă de caractere (OCR/ICR [1][2]) a datelor respective și afișarea acestora cu posibilități de editare, corectare și/sau validare de către operator înainte de a fi stocate/utilizate mai departe.

Abordarea avută în vedere a pornit de la câteva observații elementare. Completarea datelor în câmpurile formularelor se poate face computerizat, prin tipărire la imprimante diverse, cu fonturi diferite, prin dactilografiere la diferite mașini de scris sau prin scriere de mână de către diverse persoane, cu scris diferit. Pe de altă parte, din punctul de vedere al tipului de date conținute, unele câmpuri completate pot fi de tip exclusiv numeric, altele exclusiv text, iar altele mixt alfa-numeric. De asemenea, pot conține exclusiv majuscule sau minuscule în componenta de tip text, pot avea un format fix (obligatoriu) pentru caracterele completate, etc.

Astfel, o primă idee de la care s-a pornit în abordarea aleasă a constat în utilizarea unor metode și algoritmi care permit în același timp flexibilitate și generalitate în ceea ce privește partea de OCR/ICR prin folosirea unor mecanisme de învățare-clasificare și recunoaștere [3] bazate pe acumularea evolutivă, organizată, a cunoștințelor din exemple de caractere (litere și/sau cifre), ca alternativă mai adecvată decât „înghețarea” unor descrieri ale acestora prin programare (*hardcoded*). Această opțiune a fost motivată de faptul că formularele pot prezenta câmpuri variate, cu conținut diferit, care la rândul său poate fi completat în numeroase moduri, atât ca tip de scriere cât și calitate a acesteia. În acest context, alegerea unui tip adecvat de reprezentare parametrică a caracterelor constituie una din problemele cheie, ea trebuind să asigure premisele unui compromis acceptabil între puterea de generalizare și capacitatea de discriminare. A fost utilizată o reprezentare „vagă” pentru fiecare caracter, obținută pe baza porțiunii de imagine din interiorul dreptunghiului de încadrare a acestuia și codificată ca vector de N parametri. Pornind de la astfel de reprezentări în urma procesului de învățare, algoritmul de clasificare utilizat (o variantă ne-neurală a [4] cu structuri spațiale hiperparalelipipedice pentru clasificator și utilizând distanțe Manhattan în loc de euclidiene, mai eficientă atât ca și capabilitate de acoperire cât și din punct de vedere computațional) realizează acoperirea spațiului problemei cu regiuni compuse din unul sau mai multe domenii hiperparalelipipedice imbricate. Fiecare astfel de structură (regiune) specifică unei clase învățate și etichetată cu codul ASCII al caracterului alfanumeric respectiv este definită de reuniunea domeniilor hiperparalelipipedice de acoperire a clasei respective, deformate/alungite neuniform pe axe în procesele de învățare evolutivă. Structura compusă din totalitatea regiunilor astfel construite la învățare formează o bază de cunoștințe de caractere învățate. La recunoaștere, algoritmul caută clasa a cărei structură (regiune) „acoperă” reprezentarea obținută pentru caracterul țintă (de recunoscut) curent, furnizând în caz de succes codul ASCII al acestuia sau un cod ales convențional în cazul nerecunoașterii. S-a dovedit că algoritmul utilizat are o putere rezonabilă atât de generalizare cât și de discriminare în același timp, permițând o stabilizare relativ rapidă a ratei de recunoaștere. Este de subliniat totuși faptul că pentru a se atinge un nivel de performanță, stabilitate și fiabilitate convenabil este necesară instruirea sistemului pe un număr considerabil de exemple reale, mai ales pentru recunoașterea scrisului de mână.

O a doua idee importantă degajată din analiza problemei a fost aceea că, pentru fiecare tip de

formular (și condiții identice de rezoluție la scanare) să se asocieze un fișier șablon/machetă (*template*). În acesta sunt stocate informații descriptive – atribute de localizare, denumire, tip, format etc., – referitoare la fiecare câmp de interes și tipul conținutului acestuia. Aceste cunoștințe preliminare despre câmpuri și conținutul lor sunt folosite (când există) la recunoaștere, atât pentru a se diminua riscul unor confuzii, cât și pentru a se crește calitatea acesteia prin utilizarea numai a acelor baze de cunoștințe potrivite pentru tipul fiecărui câmp. Mai exact s-a avut în vedere utilizarea, atunci când este posibil, a unor baze de cunoștințe diferite specializate pentru caractere litere și respectiv cifre, incluzând aici scrisul de mână cu caractere separate. Alegerea automată a acestor baze de cunoștințe la recunoașterea anumitor câmpuri se face pe baza informației de tip asociat câmpului. În plus, pe de altă parte, rezultatul recunoașterii brute mai este supus și unor corecții automate post-recunoaștere bazate pe eventuala informație de format specific asociată câmpului respectiv, dacă aceasta există în fișierul machetă/șablon amintit mai sus.

Trebuie neapărat menționat aici că recunoașterea automată eficientă a formularelor nu poate avea loc fără o anumită „disciplină” impusă *design*-ului acestora, modului de completare, de manipulare și de achiziție a imaginilor acestora. În acest sens este de preferat ca rubricile sau chenarele în care trebuie completată informația să fie tipărite cu o culoare sensibil mai puțin intensă (de exemplu verde deschis) decât cea cu care se va realiza completarea lor (de exemplu, albastru închis sau negru). Caracterele completate, litere sau cifre trebuie să fie cât mai lizibile și cât mai rezonabil poziționate în câmpurile respective, fără a se atinge sau suprapune între ele și fără a depăși sau atinge, pe cât posibil, eventualele chenare și rubricații de separare. Formularele nu trebuie să prezinte adnotări, mângăleli și/sau alte pete. Scanarea trebuie să fie făcută cât mai îngrijit cu putință, prin setarea cât mai potrivită a parametrilor de rezoluție, luminozitate și contrast, ca și prin asigurarea unei cât mai bune reproductibilități a poziționării formularelor respective în imaginea digitală rezultată. Astfel de cerințe se înscriu în conceptul deja răspândit de „*machine readability*” pentru documente având ca scop facilitarea proceselor de analiză și segmentare automată a imaginilor, care preced recunoașterea propriu-zisă.

La nivel de arhitectură generală, pentru implementarea ideilor prezentate mai sus sistemul experimental realizat [6] conține următoarele componente funcționale:

- un subsistem de achiziție a imaginilor formularelor (prin scanare, sau de pe disc);
- un subsistem de machetare, pentru crearea și stocarea de șabloane (fișiere conținând câte o macro-descriere) specifice fiecărui tip de formular, care permite definirea interactivă, *off-line*, a câmpurilor de interes din acestea, precum și „etichetarea” lor cu atribute pe baza unor cunoștințe apriori despre tipul și caracteristicile datelor care vor fi completate, incluzând și o gramatică de descriere a acestora, bazată pe expresii regulate [5];
- un subsistem de învățare interactivă, *off-line*, dedicat instruirii pe caractere (litere sau cifre) de diverse tipuri și fonturi, preluate din câmpurile formularelor, cu crearea și actualizarea unor baze de cunoștințe corespunzătoare;
- un subsistem de recunoaștere automată a caracterelor completate în câmpurile formularului curent, utilizând informațiile despre acestea preluate dintr-un fișier machetă / șablon corespunzător precum și baza sau bazele de cunoștințe de caractere adecvate;
- un subsistem de afișare–editare–validare a rezultatelor recunoașterii pe câmpurile formularului care permite operatorului uman să compare rezultatul recunoașterii cu originalul, să facă eventuale corecții și/sau să valideze aceste rezultate;
- un subsistem supervisor de definire și configurare interactivă a contextului de lucru la un moment dat (tipul formularului – șablonul cu care se va lucra, bazele de cunoștințe ce vor fi folosite la învățare-recunoaștere etc.), comandă a funcțiilor sistemului și control automat al fluxului de lucru în contextul stabilit.

În Figura 1 este redată arhitectura și schema funcțională sintetică a întregului sistem [6].

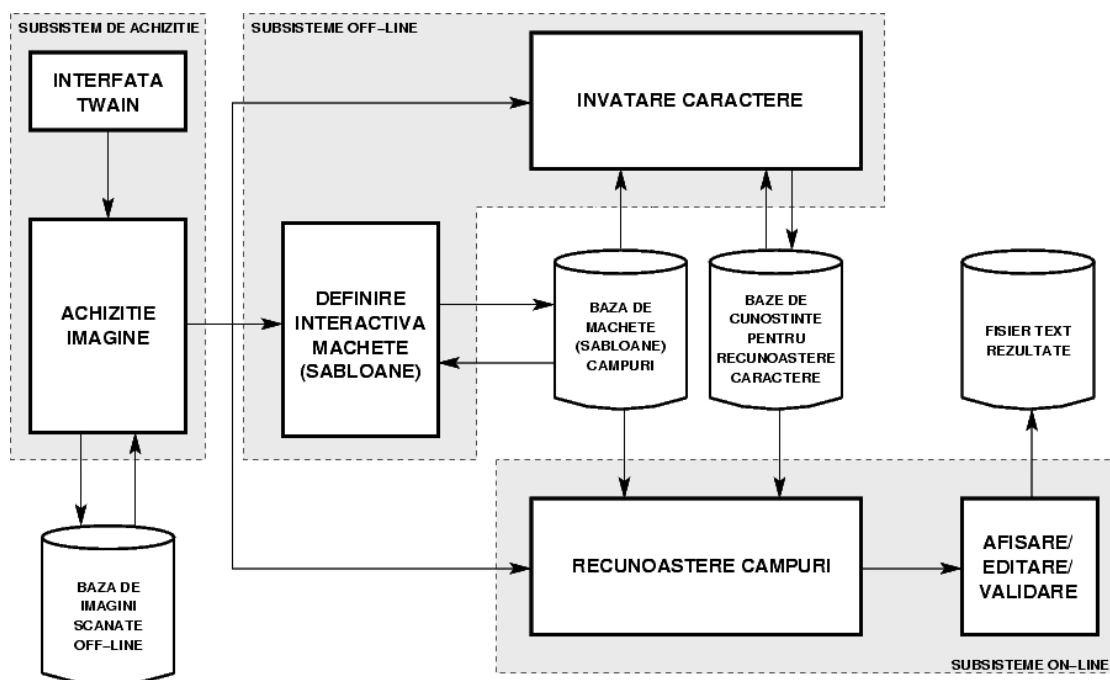


Figura 1. Arhitectura și schema funcțională sintetică a sistemului

Subsistemul de achiziție a imaginilor formularelor permite fie preluarea/digitizarea unui nou document (comandând un scanner printr-o interfață TWAIN) și salvarea acestuia, fie încărcarea de pe disc a unei imagini de formular scanat și salvat anterior. Imaginea este pusă la dispoziția subsistemelor de machetare, învățare, recunoaștere și afișare-editare-validare rezultate. Este posibilă selectarea simultană și încărcarea apoi automată, succesivă, a mai multor imagini de pe disc. Aceasta permite lucrul pe loturi (*batch*) de imagini în activitatea de introducere de date, astfel încât să se asigure o productivitate sporită.

Subsistemul de machetare oferă posibilitatea predefinirii zonelor de interes (a câmpurilor care trebuie completate) dintr-un anumit tip de formular prin reținerea coordonatelor acestora în imagine. De asemenea, permite asocierea unor atribute acestor câmpuri:

- număr de ordine câmp;
- denumire câmp;
- tip câmp (numeric / alfa / mixt / nespecificat);
- descriere codificată format câmp (dacă informația respectivă este disponibilă) printr-o gramatică de structură a câmpului respectiv, bazată pe formalismul RE (*Regular Expressions* = expresii regulate) [5], de exemplu: cu sau fără spații, exclusiv majuscule sau minuscule, specificare format fix obligatoriu etc.

Aceste atribute permit folosirea bazelor de cunoștințe cele mai potrivite precum și anumite corecții automate post-recunoaștere care să elimine unele eventuale confuzii (cum ar fi B cu 8, O cu 0, Z cu 2, l cu 1 etc.) și/sau să „sincronizeze” șirul recunoscut cu un șablon predefinit obligatoriu. Odată stabilite cele de mai sus, meta-informația astfel asociată unui anumit tip de formular este stocată pe disc într-un fișier machetă/șablon urmând a fi folosită de subsistemele de învățare și recunoaștere ori de câte ori se lucrează cu respectivul tip de document. De asemenea, ea poate fi refolosită și de subsistemul de machetare însuși pentru eventuale modificări ulterioare. Trebuie precizat faptul că, într-o astfel de machetă coordonatele imagine ale câmpurilor definite sau selectate sunt strâns legate de rezoluția de scanare, nefiind valabile decât pentru formulare de același tip, scanate la o aceeași rezoluție. Menționăm că pentru fiecare tip de câmp, pentru a se mări și mai mult gradul de încredere în preluarea automată a informației și a se diminua efortul aferent etapei de corecție-validare, mai este posibilă și asocierea câte unui dicționar, pe baza căruia să se mai realizeze o corecție automată post-recunoaștere.

Subsistemul de învățare are rolul de a acumula organizat cunoștințele necesare recunoașterii ulterioare a caracterelor alfanumerice (litere și cifre) completate în câmpurile formularelor, pornind de la exemple reale. Fiecare câmp din imaginea curentă este segmentat automat în caractere și se încearcă mai întâi o recunoaștere a acestora utilizându-se starea curentă a bazei de cunoștințe de lucru. Succesiv, pentru fiecare caracter (încă) nerecunoscut și care se dorește a fi învățat este generată reprezentarea sa parametrică, pe baza conținutului efectiv al imaginii acoperite de fereastra sa de încadrare. Setul (vectorul) de parametri obținuți este folosit de algoritmul de clasificare-învățare care are ca finalitate extinderea “experienței” bazei de cunoștințe selectate în sensul recunoașterii și a caracterului considerat. În sens larg, aceasta se poate traduce prin eventuala adăugare a unei noi instanțe pentru clasa corespunzătoare caracterului respectiv, etichetată cu codul ASCII al acestuia în baza de cunoștințe indicată. De asemenea, s-a avut în vedere și posibilitatea ștergerii/excluderii unui caracter din respectiva bază de cunoștințe la nevoie (de exemplu, în situația unei învățări greșite). Bazele de cunoștințe astfel create și/sau actualizate sunt stocate pe disc și vor fi folosite la recunoaștere și/sau pentru actualizare prin învățări ulterioare.

Subsistemul de recunoaștere analizează fiecare câmp din imaginea curentă pe baza coordonatelor definite în macheta corespunzătoare tipului de formular respectiv, realizează segmentarea acestuia în caractere și utilizează atributele asociate câmpului și bazele de cunoștințe indicate pentru a furniza textul completat (litere și/sau cifre) care a fost recunoscut, ca șir de caractere (secvență de coduri ASCII). Împreună cu acest șir recunoscut sunt furnizate mai departe subsistemului de afișare-editare-validare și informațiile obținute din macheta utilizată curent, respectiv coordonatele câmpului și celelalte atribute asociate lui, pentru toate câmpurile definite în macheta respectivă.

Subsistemul de afișare-editare-validare oferă o interfață prin care operatorul uman poate edita sau corecta rezultatele recunoașterii din fiecare câmp, având totodată în față și imaginea reală a câmpului respectiv. De asemenea poate vizualiza la dorință imaginea întregului formular. Odată acceptată corectitudinea informației preluate din formular, operatorul validează datele respective care vor fi salvate pe disc. Cu ajutorul unor interfețe de conversie potrivite, aceste date pot fi stocate/înmagazinate în formatul de bază de date specific fiecărei aplicații concrete pe care ar deservi-o sistemul de asistare a introducerii de date din formulare, sau pot fi furnizate într-un format intermediar general utilizat (cum ar fi XML), către alte aplicații specifice de prelucrare și/sau stocare a datelor respective.

Subsistemul de configurare, comandă și control (supervizor) are ca rol:

- integrarea celorlalte subsisteme implementate, într-un ansamblu funcțional unitar și coerent;
- asigurarea selecțiilor și setărilor necesare fiecăruia din acestea;
- gestionarea fluxurilor de comandă și de date implicate.

Acest subsistem furnizează practic și interfața utilizator din care pot fi comandate principalele funcțiuni ale sistemului:

- scanare imagini și salvare pe disc sau selecție imagine/imagini de lucru din cele stocate și încărcarea imaginii curente;
- definire machetă pentru un anumit tip de formular și salvare pe disc;
- selecție și încărcare machetă curentă;
- selecție și încărcare set de baze de cunoștințe curent;
- pornire sesiune de învățare din imaginea curentă;
- pornire sesiune de recunoaștere pe setul de imagini selectate, folosind macheta curentă și setul de baze de cunoștințe curente și încheiată cu afișare/editare/validare succesivă pentru fiecare imagine analizată;

- ieșire din sistem (terminare program).

În plus, subsistemul asigură și toată ordinea și logica de apelare a funcțiilor și procedurilor specifice pentru realizarea fiecăreia din aceste funcțiuni precum și managementul și controlul datelor și fluxului acestora.

Prin integrarea subsistemelor componente amintite mai sus a fost realizat un system experimental [6], care a fost utilizat pentru experimentări și testări.

3. Implementare, rezultate, observații și concluzii

Menționăm că toate implementările – pentru subsistemele componente, cât și pentru sistemul experimental – s-au realizat pentru platforme Microsoft Windows pe 32 biți, folosind limbajul de programare C.

Sintetic, rezultatele pot fi enumerate după cum urmează:

- a fost proiectat și implementat un sistem experimental pentru asistarea introducerii în calculator a informațiilor de tip text completate în câmpurile formularelor. Acesta a fost realizat prin integrarea tuturor subsistemelor descrise anterior într-un ansamblu funcțional unitar și coerent, sub o interfață grafică utilizator ergonomică, bazată pe ferestre, meniuri și *dialog box*-uri sub Windows 32 biți. Sistemul experimental implementează modul în care se realizează selecțiile, configurarea și comanda / controlul acestuia prin intermediul interfeței utilizator, permițând stabilirea contextului de lucru (ce imagine se folosește la un moment dat, ce machetă, ce baze de cunoștințe, ce moduri speciale de analiză a imaginilor, etc.), precum și lansarea funcțiilor sistemului;
- utilizându-se sistemul experimental au fost efectuate o serie de teste vizând fiecare componentă a acestuia. Au fost efectuate mai multe teste în ce privește modul de lucru cu interfața și au fost analizate rezultatele obținute pe diferite formulare;
- s-au efectuat teste care au pus în evidență faptul că algoritmul de învățare-clasificare și recunoaștere proiectat și implementat are o putere rezonabilă atât de generalizare, cât și de discriminare în același timp și permite o stabilizare relativ rapidă în cazurile „normale”. Pe de altă parte, așa cum era de așteptat, în cazuri mai „dificile”, cum este recunoașterea scrisului de mână, pentru a se atinge un nivel de performanță, stabilitate și fiabilitate convenabil este necesară instruirea sistemului pe un număr considerabil de exemple reale (multe instanțe pentru fiecare clasă).

Mai prezentăm în continuare alte câteva observații și concluzii desprinse pe parcursul experimentelor/testelor derulate.

În urma analizei rezultatelor obținute, s-au degajat câteva idei vizând necesitatea extinderii implementărilor anterioare pentru rezolvarea a două probleme:

- în cazul formularelor de același tip, scanate cu poziționări variate și/sau alegându-se setări de contrast și/sau luminozitate (praguri de binarizare) diferite etc., s-au constatat deplasări ale coordonatelor definite (prin șablon/machetă) pentru câmpuri. Deoarece funcționarea sistemului se bazează tocmai pe o corectă reproductibilitate a localizării câmpurilor în formulare de același tip, scanate în aceleași condiții de rezoluție, aceasta ducea la scăderea artificială a ratei de recunoaștere pe câmpurile respective, din cauza erorilor de segmentare.
 - pentru rezolvarea acestei categorii de probleme, s-a proiectat, implementat și adăugat o metodă prin care pentru fiecare imagine se caută un punct de referință (x_0, y_0) relevant pentru conținutul invariant acesteia, care să fie (cât mai) reproductibil de la o scanare la alta. Aceasta se realizează prin analiza matricei imaginii scanate și identificarea celei mai din stânga coloane (x_0) și respectiv celei mai de sus linii (y_0) , care conțin porțiuni consistente din conținutul invariant (tipărit) al formularului. Coordonatele câmpurilor se rețin și se utilizează ca fiind

relative la acest punct (ale cărui coordonate, x_0 și y_0 , sunt folosite ca *offset*-uri pe axele x și respectiv y);

- în cazul formularelor care prezintă marcaje / liniaturi de tip „grilă” pentru câmpurile de completat s-au observat situații în care mecanismele de segmentare a caracterelor completate în acestea nu erau suficient de eficiente.
 - ca soluție pentru acest tip de probleme s-a proiectat, implementat și adăugat o funcție de filtrare automată a liniilor orizontale (lungi) și a celor verticale (scurte) ce compun astfel de marcaje / liniaturi (de tip „grilă”).

De asemenea, așa cum am mai amintit s-a degajat pregnant necesitatea antrenării sistemului pe un număr important de exemple de caractere, mai cu seamă în cazul scrisului de mână, pentru a se putea ajunge la o anumită stabilizare a gradului de generalizare și implicit a ratei de recunoaștere automată în aceste situații.

4. Posibile dezvoltări ulterioare

Încheiem prezentând câteva posibile dezvoltări/îmbunătățiri pe care le întrezărim și care ar mai putea fi eventual aduse sistemului. Acestea pot fi grupate în următoarele trei categorii:

(a) îmbunătățirea performanțelor OCR/ICR, de exemplu prin:

- rafinarea mecanismelor de segmentare a caracterelor din zone marcate cu chenar sau linii grilă de ghidare, respectiv de filtrare mai eficientă a acestor marcaje parazite;
- extinderea mecanismelor de ajustare automată a poziției (localizării) și/sau dimensiunilor câmpurilor față de cele din șablon/machetă în funcție de eventuale încadrări puțin diferite în unele scanări;
- adăugarea unor mecanisme de corecție – validare automată pentru unele câmpuri pe baza unor dicționare atașate acestora prin șablon/machetă;
- extinderea setului de atribute de format câmp cu: câmp numeric negativ, câmp numeric cu N zecimale, câmp din care sunt relevante doar primele/ultimele N caractere;
- proiectarea și implementarea unei alte reprezentări și mai potrivite pentru caractere (mai ales pentru a asigura o generalizare mai ușoară în cazul scrisului de mână) și/sau chiar a unui algoritm alternativ de clasificare-recunoaștere complementar (eventual bazat pe memorii asociative) în acest caz;
- antrenarea intensivă și extensivă (masivă) a unor baze de cunoștințe specifice, mai ales pentru scrisul de mână;

(b) extinderea/îmbunătățirea ergonomiei interfeței utilizator de exemplu prin:

- adăugarea unui obiect de tip “*selection tree*” în interfața de corectare – editare – validare care să permită o navigare mai intuitivă printre rezultatele „OCR-izării” câmpurilor unui formular;
- adăugarea unor obiecte de control de tip *entry field*, *check box*/*radio button*, în interfața de editare atribute câmp (pentru extensiile legate de șablon/machetă: format câmp, dicționar atașat unui câmp);

(c) adaptarea funcționalității pentru aplicații specifice de exemplu prin:

- pregătirea pentru interfațarea cu – respectiv export al datelor de ieșire spre – formate specifice de baze de date, eventual printr-un format intermediar adecvat, general acceptat și utilizat (cum ar fi XML).

Potențialii utilizatori / beneficiari ai unui astfel de sistem adaptat și configurat specific pot fi: administrații financiare, poliție, evidența populației, bănci, oficii poștale, regii, agenții, operatori de telefonie, alte birouri de funcționari publici etc., unde este necesară preluarea pe scară largă a

informațiilor completate în câmpurile unor formulare tipizate, cu format fix (declarații tip, cereri tip, mandate tip, avize tip, foi de depunere / vărsământ tip, ordine de plată/încasare/schimb valutar tip etc.).

BIBLIOGRAFIE

1. **ONȚANU, D.-M.; VREJOIU, M. H.:** Sistem de recunoaștere optică a caracterelor bazat pe rețele neurale – produs program pentru recunoașterea scrisului de mână, Tema A15, Institutul Național de Cercetare-Dezvoltare în Informatică - ICI București, 1993.
2. **VREJOIU, M. H.; ONȚANU, D.-M.:** Sisteme de programe de tip OCR, PC World România, nr. 6, edit. IDG România, iunie 1995.
3. **MITCHELL, T.:** Machine Learning, McGraw-Hill, ISBN: 0070428077, March 1997.
4. **ONȚANU, D.-M.:** Learning by Evolution. A New Class of General Classifier Networks and Their Training Algorithm, Advances in Modelling & Analysis, AMSE Press, vol. 26, nr. 2, pp. 27-30, 1993.
5. **FRIEDL, J. E.:** Mastering Regular Expressions, 2nd Ed., O'Reilly, July 2002.
6. **VREJOIU, M. H.; ONȚANU, D.-M.:** Sistem de recunoaștere optică de caractere pentru citirea automată de formulare scanate. Faza a V-a. Realizare sistem experimental pentru asistarea introducerii în calculator a informațiilor de tip text completate în câmpurile formularelor, Raport de fază, proiect PN0313-0301, Institutul Național de Cercetare-Dezvoltare în Informatică - ICI București, noiembrie 2005.