

CUPRINS

1. Reprezentarea digitală a informației.....	8
1.1. Reducerea datelor.....	10
1.2. Entropia și câștigul informațional.....	10
2. Identificarea și analiza tipurilor formatelor de documente pe suport electronic.....	12
2.1. Tipuri de formate pentru documentele pe suport electronic.....	13
2.1.1 Formate text.....	13
2.1.2 Formate imagine.....	26
2.1.3 Formate audio.....	43
2.1.4 Formate video.....	48
2.1.5. Formate multimedia.....	55
2.2. Comparație între formate.....	72
3. Conținutul web.....	74
3.1. Identificarea formatelor de documente pentru conținutul web.....	74
4. Conversia documentelor din format tradițional în format electronic.....	101
4.1. Metode de digitizare.....	104
4.1.1 Cerințe în procesul de digitizare.....	104
4.1.2. Scanarea.....	104
4.1.3. Fotografieri digitală.....	105
4.1.4. Scanarea fotografiilor analogice.....	106
4.2. Arhitectura unui sistem de digitizare.....	106
4.3. Moduri de livrare a conținutului digital.....	108
4.3.1. Text sub imagine.....	108
4.3.2. Text peste imagine.....	108
4.4. Alegerea unui format.....	108
4.5. Tehnici de conversie prin Recunoașterea Optică a Caracterelor (Optical Character Recognition-OCR).....	109
4.5.1. Scurt istoric al conversiei documentelor din format tradițional prin scanare și OCR-izare.....	109
4.5.2. Starea actuală a tehnologiei OCR.....	110
4.5.3. Tehnologii curente, produse software pentru OCR-izare.....	112
4.5.4. Suportul lingvistic al produselor pentru OCR-izare.....	114

4.5.5. Particularități ale sistemelor OCR actuale.....	118
5. Elaborarea unui model conceptual pentru un sistem de informare și documentare.....	126
5.1. Depozite digitale.....	126
5.1.1. Organizații implicate în preservarea conținutului digital.....	126
5.1.2. Preocupări actuale pentru implementarea depozitelor digitale instituționale în universitățile din România.....	128
5.2. Elaborarea modelului conceptual al sistemelor de informare și documentare cu conținut tehnic. Rezultate obținute în cadrul UTBv.....	129
5.3. Implementarea unui depozit instituțional digital la partenerul Universitatea Transilvania din Brașov.....	129
5.3.1. Misiunea și obiectivele depozitului digital.....	130
5.3.2. Tehnologii pentru alegerea platformei și a software-ului.....	134
6. Concluzii.....	145
7. Note bibliografice și webografice.....	146

OBIECTIVE GENERALE ALE PROIECTULUI

Prin cercetările ce se vor întreprinde pe parcursul derulării proiectului TEMACOD, acesta își propune realizarea următoarelor **obiective generale**:

- a) Îmbunătățirea, coordonarea și eficientizarea proceselor de digitizare a resurselor informaționale din domeniul tehnic din universități și alte tipuri de organizații;
- b) Sporirea numărului de resurse informaționale tehnice reprezentative digitizate, diversificarea și conservarea acestora;
- c) Creșterea gradului de accesibilitate a publicului utilizator de Internet la resursele informaționale tehnice.

Pentru atingerea acestor obiective proiectul **TEMACOD** are în vedere și realizarea următoarelor **obiective specifice** în cadrul fiecărui obiectiv general:

Obiectivul general (a)

1. Identificarea sistemelor de informare și documentare care vehiculează documente tehnice și coordonarea modalităților de digitizare a resurselor informaționale în organizații;
2. Soluționarea și coordonarea aspectelor logistice și documentare pe care le presupune digitizarea pentru înlăturarea paralelismelor și suprapunerilor, pe fiecare arie tematică în parte;
3. Crearea unor mecanisme de coordonare, gestionare și întreținere a procesului de digitizare, care să permită o arhitectură coerentă a întregului sistem, la nivel de proiect.

Obiectivul general (b)

1. Identificarea unui segment documentar reprezentativ de resurse informaționale, pe fiecare arie tematică;
2. Dezvoltarea unui sistem performant și eficient în domeniul tehnologiei informației, care să permită crearea unei platforme unice și unitare de acces la resursele digitizate, precum și conservarea materialului digitizat.

Obiectiv general (c)

1. Aplicarea reglementărilor legale referitoare la drepturile de autor și la proprietatea intelectuală a resurselor informaționale;
2. Instituirea unui sistem unic și unitar de catalogare și/sau clasificare și descriere a resurselor informaționale, pe fiecare arie tematică în parte.

Programul național de informatizare a societății românești, cu un sens mai larg de integrare a culturii și științei naționale în circuitele mondiale, generează în mod natural și logic un nou tip de acces la valorile tehnice și științifice din țară și din lume, un nou tip de relații între marea categorie a celor dornici de a ști și grupul restrâns ca dimensiuni al celor care oferă informații: instituțiile educaționale, de cercetare, de cultură și artă.

În cadrul acestui program, un rol important revine procesului de informatizare a bibliotecilor, care a trecut de la faza de înregistrare a datelor în cataloage, la cea de dezvoltarea de sisteme integrate de bibliotecă. În prezent ne aflăm în etapa de constituire și dezvoltare a colecțiilor de date digitale în biblioteci, care va permite ulterior interconectarea bibliotecilor digitale în rețele, bazate pe infrastructuri de comunicații.

Proiectul propus în cadrul Programului 4, „Parteneriate în domeniile prioritare”: *Tehnici pentru managementul conținutului digital-TEMACOD* are ca obiective principale: dezvoltarea de biblioteci digitale, bazate-Web, la nivelul instituțiilor de învățământ superior implicate în contract, interconectarea lor și accesarea rapidă via Internet, printr-o interfață prietenoasă și eficientă.

Tema abordată este de mare actualitate în întreaga lume, se pune problema construirii unei Biblioteci digitale la nivel mondial, accesibilă prin Internet, gratuită și în format multilingv, care să conțină materiale semnificative.

În constituirea acestei biblioteci digitale mondiale, un rol important va reveni mediului academic, bogat în resurse educaționale, de cercetare și care are deja o mare parte din fondul infodocumentar digitizat.

Spre deosebire de bibliotecile tradiționale, o bibliotecă digitală are colecțiile stocate în format digital și accesibile prin calculatoare conectate local sau la distanță, prin rețele de calculatoare. De asemenea, ea oferă un sistem performant de regăsire a informației.

Ca rezultat al dezvoltării Internetului și a potențialului său de căutare, bibliotecile digitale cum ar fi Biblioteca Europeană și Biblioteca Congresului sunt acum dezvoltate într-un mediu bazat-Web. Bibliotecile publice, cele școlare și academice sunt de asemenea capabile să dezvolte website-uri care să permită preluarea de informație științifică digitizată, din categoria cărților scrise, a celor audio sau video, schimbând fundamental noțiunea de resursă educațională.

Multe biblioteci academice sunt activ implicate în constituirea depozitelor de cărți, buletine științifice, teze de doctorat și alte lucrări care pot fi digitizate sau au fost de la început preluate în format digital. Aceste depozite sunt accesibile în rețeaua universității, dar și în afara ei, cu anumite restricții, ce se referă în special la protejarea rezultatelor cercetării și care sunt în concordanță cu politica adoptată de instituție privind accesul liber.

OBIECTIVELE ETAPEI DE EXECUȚIE

ACTIVITATE: Analiza formatelor documentelor tehnice pe suport electronic

Obiectiv: Identificarea tipurilor de formate în care se regăsesc documentele pe suport electronic

ACTIVITATE: Analiza normelor de conversie a conținutului cuprins în documentele tehnice de la formatul tradițional la formatul electronic prin tehnici de digitizare

Obiectiv: Identificarea celor mai noi tehnici de trecere a documentelor din format tradițional în format electronic. Definirea instrumentelor pentru realizarea depozitelor de conținut digital.

ACTIVITATE: Analiza conținutului web; analiza formatelor documentelor tehnice regăsite pe web

Obiectiv: Realizarea unui studiu asupra tipurilor de conținut web și a formatelor documentelor tehnice

ACTIVITATE: Elaborarea modelului conceptual al sistemelor de informare și documentare cu conținut tehnic.

Obiectiv: Elaborarea principiilor de bază și a metodologiei identificării și creării sistemelor infodocumentare pentru documente cu conținut tehnic.

REZUMATUL ETAPEI

Proiectul de cercetare TEHNICI PENTRU MANAGEMENTUL CONȚINUTULUI DIGITAL și-a propus în etapa 2, conform Planului de activități, atingerea obiectivelor prin finalizarea a patru activități de bază și anume: analiza formatelor documentelor tehnice pe suport electronic, analiza normelor de conversie a conținutului cuprins în documentele tehnice de la formatul tradițional la formatul electronic prin tehnici de digitizare, analiza conținutului web-analiza formatelor documentelor tehnice regăsite pe web, elaborarea modelului conceptual al sistemelor de informare și documentare cu conținut tehnic.

Activitățile de cercetare efectuate au avut ca scop identificarea și analiza tipurilor de formate regăsite pe suport electronic fie online, fie pe unități fixe sau mobile de memorie, care pot avea în componența lor diferite obiecte informaționale. Deasemenea s-a analizat conținutul web în care informația se găsește sub formă nestructurată, dar și structura și importanța sistemelor de gestiune de baze de date în care regăsim informație structurată.

Partea aplicativă a activității de cercetare a constat în definirea unui model conceptual pentru managementul informației tehnice în sisteme infodocumentare, implementat la Universitatea Transilvania din Brașov. Acesta va constitui un depozit digital care va cuprinde documente specifice activității de cercetare și educaționale din universitate.

Metoda de cercetare a avut la bază un studiu teoretic asupra tipurilor și analizei tipurilor de formate a documentelor cel mai des utilizate în mediul informatizat. Deasemenea a avut loc și o cercetare practică desfășurată cu preponderență la partenerul Universitatea Transilvania din Brașov în vederea definirii și implementării unui depozit digital instituțional.

Pe perioada desfășurării etapei 2 au avut loc întâlniri cu partenerii, workshop-uri, dezbateri ale aspectelor identificate spre rezolvare.

Raportul de cercetare pentru această etapă poate fi considerat ca fiind o sinteză din punct de vedere teoretic a principalelor activități ale etapei, cât și inițierea punerii în aplicație a concluziilor teoretice.

Analiza formatelor documentelor tehnice pe suport electronic. Munca de alegere a unui anumit format este îngreunată de multitudinea de opțiuni disponibile. S-au studiat mai multe posibilități, luând în considerare avantajele și dezavantajele fiecăruia și având mereu în vedere scopul utilizării fișierului. S-au identificat formate pentru principalele tipuri de obiecte informaționale conținute: text, imagini, audio, video, multimedia și s-au evidențiat principalele caracteristici pentru care se optează pentru un format sau altul.

Analiza normelor de conversie a conținutului cuprins în documentele tehnice de la formatul tradițional la formatul electronic prin tehnici de digitizare. În acest capitol s-a făcut o

scurtă prezentare a diferitelor modalități în care se poate realiza achiziția imaginilor în formă digitală, plecând de la o carte. Deasemenea s-a definit o arhitectură a unui sistem de digitizare și s-au dezvoltat tehnicile de conversie prin Recunoașterea Optică a Caracterelor (Optical Character Recognition-OCR). Pe lângă un scurt istoric al conversiei documentelor din format tradițional prin scanare și OCR-izare, am prezentat tehnologiile curente și produse software pentru conversia în format digital. S-a evidențiat și suportul lingvistic al produselor OCR cât și alte particularități.

Analiza conținutului web; analiza formatelor documentelor tehnice regăsite pe web. S-a analizat tipul de conținut web și creatorii de conținut online. Cele mai frecvent utilizate formate de realizare a paginilor web au constituit baza de plecare în analiza tipurilor de documente care circulă în rețeaua Internet sau altor tipuri de rețele.

Elaborarea modelului conceptual al sistemelor de informare și documentare cu conținut tehnic. Cercetare din etapa 2 se încheie cu elaborarea unui model de management al informației în structuri infodocumentare cu conținut preponderent tehnic. S-a pornit de la definirea depozitelor digitale și a organizațiilor care dezvoltă crearea și prezervarea pe termen lung a arhivelor digitale.

S-au evidențiat și preocupările actuale pentru implementarea depozitelor digitale în universitățile din România, activitate ce este la început de drum. Pentru realizarea depozitului digital la Universitatea Transilvania din Brașov s-a folosit platforma de cercetare interdisciplinară ASPECKT.

S-au evidențiat tehnologiile pentru alegerea platformelor și a software-urilor pentru realizarea Depozitului pilot. Pagina web a depozitului digital ASPECKT-Dspace se poate accesa la adresa <http://aspeckt.unitbv.ro/dspace>.

Aparatul critic al raportului de cercetare este constituit din bibliografie selectivă și webgrafie selectivă.

1. Reprezentarea digitală a informației

Teoria informației răspunde la două întrebări fundamentale din teoria comunicației: care este cea mai bună compresie a datelor (răspuns: entropia H) și care este cea mai bună rată de comunicație pentru transmiterea datelor (răspuns: capacitatea canalului C). Din acest motiv unii consideră teoria informației ca fiind o submulțime a teoriei comunicației. Într-adevăr, ea are contribuții fundamentale în fizica statistică (termodinamică), știința calculatoarelor (complexitatea Kolmogorov sau complexitatea algoritmilor), inferențe statistice și la probabilitate și statistică.

La începutul anilor 1940 se credea că creșterea ratei de transmisie a informației peste un canal de comunicație crește probabilitatea erorii. Shannon a surprins comunitatea științifică dovedind că această afirmație nu este adevărată atunci când rata de comunicație este sub capacitatea canalului. Capacitatea canalului poate fi calculată simplu din caracteristicile de zgomot ale acestuia. Shannon argumenta că procese aleatoare ca muzica și vocea au o complexitate ireductibilă sub care semnalul nu poate fi comprimat. Acest lucru l-a numit el Entropie, prin asociere cu utilizarea cuvântului în termodinamică, și justifică faptul că dacă entropia sursei este sub capacitatea canalului comunicația asimptotic fără erori poate fi atinsă.

Astăzi teoria informației reprezintă punctul extrem al mulțimii tuturor schemelor de comunicație posibile. Limita compresiei de date este o extremă a setului de idei de comunicație. Toate metodele de compresie necesită descrierea datelor cel puțin la nivelul acestui minim. La cealaltă extremă este maximul transmisiei date cunoscut ca și capacitatea canalului. Astfel toate metodele de modulație și compresie a datelor se situează între aceste limite.

Teoria informației deasemenea sugerează mijloace de a obține aceste limite extreme ale comunicației. Totuși aceste scheme teoretice optime de comunicație, se dovedesc a fi nefezabile computațional. Aceasta este mai curând datorită fezabilității computaționale a schemelor de modulare și demodulare simple pe care le folosim și nu datorită codificării aleatoare și regulii de decodificare „nearest neighbor” propuse de demonstrația lui Shannon a teoremei capacității canalului. Progresele din domeniul circuitelor integrate și proiectării codului ne-au permis obținem câteva din câștigurile sugerate de teoria lui Shannon. Un bun exemplu de aplicare a acestor idei a teoriei informației este utilizarea codurilor corectoare de erori pe compact discuri.

Cercetările moderne privind aspecte ale teoriei informației în comunicație sau concentrat pe teoria informației în rețele: teoria ratelor de comunicație simultană de la mai mulți emițători la mai mulți receptori într-o rețea. Unele câștiguri de rată de comunicație între emițător și receptor sunt neașteptate și toate au o simplitate matematică certă. O teorie unificatoare rămâne totuși a fi găsită.

Știința calculatoarelor (complexitatea Kolmogorov). Kolmogorov, Chaitin și Solomonoff au avansat ideea că, complexitatea unui șir de date poate fi definită ca lungimea celui mai scurt program binar care generează acel șir de caractere. Deci complexitatea este dată de lungimea descrierii minime. Această definiție a complexității este universală, independentă de calculator și de o importanță fundamentală. Complexitatea Kolmogorov pune bazele teoriei descriptive a complexității. Din fericire, complexitatea Kolmogorov este aproximativ egală cu entropia Shannon H dacă secvența este aleasă aleator dintr-o distribuție care are entropia H . Așadar legătura între teoria informației și complexitatea Kolmogorov este perfectă. Într-adevăr considerăm complexitatea Kolmogorov a fi mai fundamentală decât entropia Shannon. Ea este limita compresiei datelor și ne duce spre o procedură consistentă logic pentru inferență.

Există aici o plăcută relație de complementaritate între complexitatea algoritmică și complexitatea computațională. Putem privi complexitatea computațională (complexitatea în timp) și complexitatea Kolmogorov (lungimea programului sau complexitatea descriptivă) ca două axe corespunzând timpului de rulare a programului și lungimii programului. Complexitatea Kolmogorov se concentrează pe minimizări de-a lungul celei de-a doua axe iar complexitatea computațională se focalizează pe minimizare de-a lungul primei axe. S-au făcut puține încercări pentru minimizarea simultană pe ambele axe.

Fizică (Termodinamică). Mecanica statistică este locul de naștere al entropiei și a celei de-a doua legi a termodinamicii. Entropia întotdeauna crește. Pe lângă alte lucruri a doua lege ne permite să rejectăm orice pretenție de „perpetuum mobile”.

Matematică (Teoria probabilității și statistică). Mărimile fundamentale ale teoriei informației – entropie relativă și informație mutuală – sunt definite ca funcționale ale distribuției probabilitice. În schimb, ele caracterizează comportamentul unor secvențe lungi de variabile aleatoare și ne permit să estimăm probabilitatea unor evenimente rare și să găsim cel mai bun exponent al erorii în testarea ipotezelor.

Filozofia științei (Briciul lui Occam). William of Occam a spus „Causes shall not be multiplied beyond necessity” sau parafrazându-l „cea mai simplă explicație este cea mai bună”. Solomonoff și ulterior Chaitin, au argumentat că se obține o procedură de predicție universal bună dacă se ia o combinație ponderată a tuturor programelor care explică datele și se observă ce tipăresc ele în continuare. Mai mult, această inferență va funcționa în multe probleme care nu sunt rezolvate prin statistică. De exemplu această procedură va prezice în final următorii digiți ai numărului π . Când această procedură este aplicată la aruncarea unei monede care generează cap cu probabilitate 0.7, acest lucru va fi de asemenea inferat. Când este aplicată la bursă procedura trebuie să găsească legile bursei și să le extrapoleze optimal. În principiu, așa o procedură ar fi găsit legea lui Newton din fizică. Bineînțeles, asemenea inferență este impracticabilă, deoarece a

elimina toate programele care nu reușesc să genereze datele existente durează exagerat de mult. Am putea prezice ce se întâmplă mâine pentru sute de ani de acum înainte.

Economie (Investiții). Investiții repetate într-o piață staționară duce la creșteri exponențiale ale bunăstării. Rata de creștere a bunăstării este duală ratei entropiei pieței. Paralela între teoria investiției optimale în burse și teoria informației este evidentă. Se poate dezvolta o teorie a investițiilor pentru a explora această dualitate.

Calcul versus Comunicație. Pe măsură ce construim calculatoare mai mari din componente mai mici atingem atât o limită de calcul cât și o limită a comunicării. Calculul este limitat de comunicare și comunicarea limitată de calcul. Acestea devin interdependente și deci toate dezvoltările din teoria comunicației prin teoria informației ar trebui să aibă un impact direct asupra teoriei computaționale.

1.1.Reducerea datelor

Selecția unei submulțimi de trăsături caracteristice este definită ca procesul de selecție a acelei submulțimi de trăsături d dintr-un set de dimensiune mult mai mare D care maximizează performanțele de clasificare peste toate submulțimile posibile. Căutarea după o astfel de submulțime este o problemă foarte dificilă. Spațiul de căutare poate fi foarte mare.

1. 2. Entropia și câștigul informațional

Câștigul informațional și entropia sunt funcții ale distribuției probabilistice care susțin procesul de comunicare. Entropia este o măsură a incertitudinii unei variabile aleatoare. Dându-se o colecție S de n eșantioane grupate în c concepte țintă (clase), entropia lui S relativă la clasificare este (1):

$$Entropy(S) = - \sum_{i=1}^c p_i \log_2(p_i) \quad (0)$$

unde p_i este procentajul din S care aparține clasei i .

Pe baza entropiei este definită o măsură a gradului de eficiență în selecția trăsăturilor. Această măsură este numită *Câștigul informațional* și reprezintă de fapt reducerea în entropie, cauzată de gruparea eșantioanelor în acord cu un atribut. Mai precis, câștigul informațional pentru un atribut relativ la o mulțime de eșantioane S este definit ca:

$$Gain(S, A) \equiv Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v) \quad (0)$$

unde $Values(A)$ este mulțimea de valori posibile pentru atributul A și S_v este submulțimea din S pentru care atributul A are valoarea v .

Utilizând ecuațiile, pentru fiecare trăsătură se calculează câștigul informațional obținut dacă mulțimea este împărțită utilizând acea trăsătură. Se obțin valori între 0 și 1 fiind apropiate de 1 dacă trăsătura împarte mulțimea originală în două submulțimi cu dimensiuni apropiate. Pentru selectarea trăsăturilor relevante am utilizat diferite praguri. Dacă câștigul informațional obținut pentru o trăsătură depășește pragul acea trăsătură va fi selectată, altfel nu va fi selectată.

Forman a arătat că câștigul informațional eșuează în cazul problemelor reale de clasificare când se lucrează cu multe trăsături.

2. Identificarea și analiza tipurilor formatelor de documente pe suport electronic

Un format de fișier reprezintă o metodă de organizare a informației în secvențe de biți pentru ca aceasta să poată fi păstrată într-un calculator. Există diferite formate de fișiere pentru documente, imagini, sunete, videoclipuri, toate acestea fiind asociate cu anumite extensii. Câteva exemple de asociații: documente cu .doc, imagini cu .JPEG, .html pentru pagini.

Munca de alegere a unui anumit format este îngreunată de multitudinea de opțiuni disponibile. Înainte de a trece la fapte, trebuie studiate mai multe posibilități, luând în considerare avantajele și dezavantajele fiecăruia și având mereu în vedere scopul utilizării fișierului.

Când considerăm toate aceste aspecte putem avea în minte următoarea listă:

- Este formatul patentat sau controlat de organizații publice de standardizare?
- Au fost făcute publice specificațiile formatului?
- Este vechimea formatului un risc?
- Este posibil un upgrade al software-ului? Dacă da, poate versiunea nouă să deschidă fișiere mai vechi?
- Este posibil ca software-ul să fie retras de pe piață?
- Este posibil ca formatul să nu mai fie folosit și, ca urmare, să dispară software-ul prin care fișierele puteau fi editate?
- Este permisă extragerea datelor dintr-un anumit format pentru a putea fi folosite mai departe sau indexate, sau este permisă doar vizualizarea acestora?
- Formatul salvează datele la calitatea dorită? Vor fi datele degradate prin salvarea în acest format?
- Sunt compresate datele în acest format? (În general compresia face ca fișierele să fie mai puțin tolerante la erori în transmiterea de date și pot fi degradate mai ușor când sunt păstrate pe diferite medii de stocare a datelor)
- Este formatul ales un standard acceptat?

2.1. Tipuri de formate

2.1.1. Formate text

DOC - fișier document

Formatul de fișier DOC este cel mai frecvent utilizat de Microsoft Word. Formatul a evoluat de-a lungul timpului astfel încât există defapt mai multe formate ce se regăsesc sub aceeași titulatură : doc.

DOC (abreviere de la *document*) este extensia numelui de fișier al documentelor produse de diverse procesoare de text în baza formatului impus de Microsoft Word versiunile 1.0-2003. Este un tip de fișier binar, care pe lângă textul propriu-zis stochează și datele legate de formatarea acestuia, de asemenea și elemente non-textuale precum imagini, grafică sau diagrame.

Este un format închis, standardul oficial pentru acest format fiind deținut exclusiv de către compania Microsoft și nefiind făcut public niciodată. Cu toate acestea, prin metode de inginerie inversă, s-a reușit reconstituirea formatului de către alți producători de software. Astăzi un număr mare de programe non-Microsoft (în special Open Source) pot să deschidă și să editeze fișiere de tip .DOC fără probleme. Printre acestea se numără Open Office și AbiWord.

Formatul Doc diferă de la o versiune la alta a Microsoft Office Word, astfel versiunile Word până la 97 au folosit un alt format DOC față de versiunile Word dintre 97 și 2003. Versiunile mai vechi sunt incapabile să citească, să primească sau să printeze corect fișierele .doc generate de noile versiuni.

Principalele operațiuni întreprinse în fișierele document:

- editarea și formatarea textelor (folosind fonturi diferite, stiluri și dimensiuni diferite de font, tipuri de subliniere a caracterelor, culoarea caracterelor etc.);
- inserare imagini, tabele, grafice, diagrame;
- așezare în pagină sau formatare (dimensiunile paginii, orientarea ei, alinierea textului în pagină etc.);
- printare;
- căutare elemente specifice (cuvinte, fraze, etc.) în cadrul textului.

Dezavantajul formatului este acela că, fiind un format închis, există încă incompatibilități de manevrare a documentelor la transferul lor între diferite procesoare de text, ceea ce poate duce la pierderea de informații de formatare și astfel, poate altera calitatea fișierului document.

Fișierele .doc sunt utilizate ca și formate-sursă pentru unele fișiere eBook. Cu ajutorul unor programe convertoare se pot converti chiar fișierele .doc direct într-un anumit format eBook. Totuși ele nu sunt folosite foarte mult ca și formate eBook deoarece folosind Word

pentru a citi documentul este foarte ușor de a-l modifica din greșeală, iar „scrolling”(derularea paginii) nu este cea mai bună metodă de a citi o carte. Unele versiuni ale Word-ului au, totuși, un „review mode” ce poate asigura o citire mai ușoară a acestor documente.

Deși este considerat învechit și a fost înlocuit inclusiv de către creatorii săi, fișierul document DOC rămâne, pentru moment, cel mai utilizat format pentru procesoare de text²⁾ cu atât mai mult cu cât acesta „oferă utilizatorului flexibilitatea necesară de a converti aproape orice tip de fișier într-un document MS-Word”.

Programele care deschid fișiere document în format DOC sunt

- în Mac OS: Microsoft Word, Apple Pages, Apple AppleWorks, OpenOffice.org Writer, Nuance OmniPage Pro X;
- în Windows: Microsoft Word, Microsoft Word Viewer, OpenOffice.org Writer, Nuance OmniPage Professional 17;
- în Linux: OpenOffice.org Writer.

Formatul de fișier DOC mai este folosit și de către WordPad Document, însă acesta nu este utilizat pe o scară atât de largă ca și Microsoft Word Document.

Începând cu Office 2007, Microsoft a înlocuit standardul DOC cu un nou standard, de data aceasta deschis (și publicat), bazat pe XML și deci folosind fișiere text în locul celor binare : .DOCX

DOCX - fișier document în Office Word 2007

DOCX este o variantă îmbunătățită de format de fișier document introdus de MS Word 2007 (parte din MS Office 2007) ce „se bazează pe tehnologiile standard XML și ZIP”.

Diferența dintre cele două formate (DOC și DOCX) constă în felul cum a fost creat fiecare și nu în conținut, astfel că, spre ex., un proces verbal, o scrisoare comercială, o invitație la o aniversare sau orice alt document ce poate fi editat în Word poate fi salvat în oricare format de fișier document.

Spre deosebire de formatul doc, DOCX este un format deschis, adică „specificațiile formatului sunt disponibile prin intermediul unei licențe care nu prevede plata unor drepturi de autor”. Prin urmare, formatul permite utilizatorilor să facă schimb de documente chiar dacă nu folosesc același procesor de text, fără riscul de a altera în vreun fel conținutul lor, ei putând deschide, edita și salva fișierele .docx.

De asemenea, „.docx nu este un fișier binar, deci nu necesită un program compatibil pentru a-l deschide, și este de dimensiuni mult mai mici decât echivalentul lui în format doc.

Un fișier .docx deși pare a fi un singur fișier, este de fapt un pachet de fișiere arhivate împreună.” Pentru a-l deschide este nevoie de o aplicație software capabilă să dezarchiveze fișierele .docx (WinZip, PKWare Unzip, și Stuffit Expander), „capabilă să înțeleagă și modul de

lucru folosit de .docx pentru a gestiona diferitele componente ale fișierului : documentul, schema, proprietățile, obiectele incluse în document etc.”

Datorită dimensiunilor sale reduse și a compatibilității cu alte procesoare de text non-Microsoft, acest tip de format de document se pretează foarte bine la trimiterea fișierelor prin poșta electronică, prin rețele sau prin Web, el putând fi deschis și citit de aproape oricine îl primește.

Totuși, chiar dacă formatul de fișier .docx este un format de fișier text ce poate fi citit cu ușurință prin dezarhivarea pachetului și deschiderea folder-ului conținând document.xml, deschiderea fișierului în acest fel alterează formatările de font și de paragraf, precum și alte obiecte integrate în fișier, astfel că se preferă folosirea unui convertor docx, care ușurează salvarea documentului dintr-un format în altul cu o configurare minimă, păstrând caracteristicile inițiale complete ale documentului. Procesul de conversie de la un format la altul implică, în esență, extragerea de text și orice obiecte înglobate, precum și stilul de formatare, apoi re-salvarea acestora în alt format.

Cele mai folosite convertoare pentru fișiere .docx sunt Zamzar, Docx Converter sau Docx2Doc, aplicații web ce fac online conversia dintr-un format în altul, nefiind nevoie de instalarea vreunui soft.

Programele care deschid fișiere document în format DOCX sunt:

- în Mac OS: Microsoft Word 2008, Apple Pages, Penergy docXConverter;
- în Windows: Microsoft Word 2007, Microsoft Word with Compatibility Pack, Panergy docXConverter, OxygenOffice Professional, OpenOffice.org with Odf-Converter, NativeWinds Docx2Rtf;
- în Linux: OxygenOffice Professional, OpenOffice.org with Odf-Converter.

DOT- fișier document Word templates

DOT este un fișier document predefinit, folosit ca șablon pentru realizarea de mai multe documente cu formate similare, cum ar fi rapoarte de afaceri, plicuri de scrisori, cereri, formulare etc.¹⁴⁾

Corporația Microsoft este proprietarul fișierului cu extensia .dot folosit la editare de texte în Word. „Aceste fișiere cuprind caracteristici de formatare standard ca: setările paginilor (mărime, orientare, margini, număr de coloane ș.a.m.d.), stiluri de fonturi, culori de text, dar și bare de instrumente (toolbars) particularizate sau macrocomenzi pentru a completa textele predefinite, ca de ex.: un nume de companie sau o adresă într-un document realizat cu ajutorul fișierului template (fișierul-șablon).

Șabloanele DOT sunt utile în special la realizarea de documente de același tip, cu aceleași formătări ori caracteristici de editare, sau la crearea personalizată a template-urilor pentru o varietate de documente ca: scrisori, faxuri, plicuri.”

De exemplu, un template pentru o scrisoare de afaceri va avea un antet (letterhead) pe prima pagină, numerotarea paginilor și eventual un câmp de generare automată a datei. În plus, pentru a personaliza stilurile pentru diferitele părți ale scrisorii (zona adresei, salutul, cuprinsul scrisorii, semnătura etc.) template-urile pot cuprinde câmpuri de introducere a textului indicându-se unde merg puse aceste elemente.

Macrocomenzile încorporate pot fi folosite pentru a completa informații în mod automat de fiecare dată când este creat un document nou folosind șablonul DOT salvat. Acestea grupează o succesiune de acțiuni ce se doresc executate direct, prin activarea unui buton, a unei combinații de taste sau eventual a unei opțiuni de meniu. Cea mai simplă modalitate de creare a unui macro este înregistrarea acțiunilor realizate, după care se stabilește modul său de execuție, optând pentru una din asocierile anterior enumerate. Astfel, grupul de acțiuni înregistrat se va repeta ori de câte ori e nevoie, cu un efort minim.

Șablonul DOT poate fi realizat prin crearea unui document în Microsoft Word și salvarea lui ca template ce poate fi deschis pe viitor folosind caseta de dialog a noului document. În plus, Microsoft Word este prevăzut cu un număr din cele mai folosite șabloane pentru scrisori, rapoarte, site-uri web și publicații.

Șablonul normal DOT este cadrul uzual pe care se construiește un document, el setează marginile și fontul Times New Roman de dimensiune 10. Pentru creșterea eficienței în crearea documentelor, se pot folosi și alte șabloane, predefinite sau definite de către utilizator.

Crearea unui șablon se face ca un document normal, acordând atenție deosebită atributelor de formatare și introducând informațiile comune tuturor documentelor ce vor fi create pe baza aceluși șablon. Fișierul creat se va salva cu tipul Template și i se va atribui extensia .dot

Deși fișierele .dot însele nu reprezintă un risc asupra securității informatice, unele pot cuprinde macrocomenzi încorporate, acestea fiind o potențială sursă de viruși informatici, astfel încât ar trebui evitate fișierele .dot conținând macro-uri din surse necunoscute.

Programele care deschid fișiere document în format DOT sunt:

- în Mac OS: Microsoft Word și NeoOffice Writer
- în Windows: Microsoft Word ¹⁴⁾

Dezavantajul formatului este acela că fiind un fișier binar și închis (ca și .dot) are dimensiuni mari și întâmpină greutatea la transferul documentelor între diverse aplicații software.

DOTX : Word 2007 Document template

DOTX este formatul noului editor de text Microsoft Word 2007 (parte din Microsoft Office 2007); el este un fișier template și „o combinație între arhitectura XML și arhivarea ZIP pentru reducerea dimensiunii.”

Prin adoptarea formatului de fișier bazat pe XML noul format template asigură o transmitere ușoară a datelor între diferite aplicații, platforme și Internet browser-i fără ca documentele astfel transmise să sufere modificări de orice fel.

Alte avantaje ale formatului Office XML:

- **Fișiere compacte.** Fișierele sunt comprimate în mod automat și pot fi până la 75 la sută mai mici, în unele cazuri. Formatele Office XML utilizează tehnologia de comprimare ZIP pentru a stoca documente, oferind economii de cost prin reducerea spațiului de disc necesar pentru stocarea fișierelor și reducând lățimea de bază necesară pentru a trimite fișierele prin poșta electronică, prin rețele și Internet. Când se deschide un fișier, acesta este dezarhivat automat, când se salvează un fișier, acesta este automat arhivat la loc;
- **Recuperarea îmbunătățită a fișierelor deteriorate.** Fișierele sunt structurate modular, menținând diferitele componente de date din fișier separate între ele. Astfel, se permite deschiderea fișierelor, chiar dacă o componentă din fișier (de ex. o diagramă sau un tabel) este deteriorat;
- **Detectare mai ușoară a documentelor ce conțin macrocomenzi.** Fișierele salvate utilizând sufixul implicit „x” (.docx, .dotx, .pptx etc.) nu pot conține macrocomenzi VBA (Visual Basic for Applications) sau controale ActiveX și astfel nu implică riscurile de securitate asociate cu aceste tipuri de cod încorporat. Doar fișierele care au extensia numelui terminată în „m” (.docm, .dotm, .xlsm etc.) pot conține macrocomenzi VBA și controale ActiveX, care sunt stocate într-o secțiune distinctă din fișier. Astfel, extensiile de nume de fișier facilitează distingerea fișierelor care conțin macrocomenzi de cele care nu conțin, ușurând identificarea de către software-ul antivirus a fișierelor ce conțin cod cu potențial dăunător;

Față de un fișier document obișnuit (.docx) la deschiderea unui fișier .dotx se deschide defapt o copie a template-ului, ceea ce protejează șablonul de modificări accidentale.

Fișierul cu extensia .dotx nu are macrocomenzi și este folosit la depozitarea documentelor electronice, cum sunt: memo-uri, rapoarte, cărți, foi de calcul, diagrame etc.

Acest format aduce îmbunătățiri privind gestionarea informației și a fișierelor, recuperarea datelor și interoperabilitatea sistemelor din rețea :

- Orice aplicație software bazată pe XML poate avea acces și lucra cu fișierele document în acest nou format, chiar dacă nu face parte din sistemul Microsoft Office sau, mai mult, chiar dacă nu este un produs Microsoft.
- Utilizatorii pot, de asemenea, folosi transformările standard pentru a elimina sau adăuga date în aceste fișiere.
- În plus, problemele legate de securitatea informatică sunt drastic reduse deoarece informațiile sunt depozitate pe suport XML, care este în esență un simplu text. Astfel, datele pot trece de programele antivirus fără nici un impediment.

Programele care deschid fișiere .dotx sunt:

- în Mac OS: Microsoft Word 2007 sau versiunile precedente cu suport Open XML
- în Windows: Microsoft Word 2007 sau versiunile precedente cu suport Open XML; OpenOffice.org Writer; Adobe LiveCycle Designer

RTF- Rich Text Format

Creat de Microsoft în 1987, ca parte a Microsoft Word 3.0 pentru Macintosh, RTF reprezintă formatul de fișier folosit pentru transferul de documente; este un format universal, adică „poate fi citit și scris de aproape toate procesoarele de text și chiar de cele mai multe editoare de text”. Acest tip de format este mult mai ușor de generat decât PDF sau PostScript, și mai prietenos pentru procesoarele de text decât HTML. Astfel că multe programe de procesoare de text permit salvarea unui document în RTF selectând opțiunea „Save As”

RTF este o aplicație ce permite utilizatorului să editeze conținutul unui document. El este standardizat pentru crearea de fișiere de text formatat. Practic „este un format de fișier ce codifică graficele și textele formatate pentru a permite transferul ușor între diferite aplicații și sisteme de operare.”

Față de fișierele de text simple, de bază, un fișier .rtf conține informații referitoare la stilul textului, mărime și culoare. RTF folosește setul de caractere ANSI (American National Standards Institute), PC-8, Macintosh sau pentru calculatoare IBM pentru a controla descrierea și formatarea documentelor. Formatul suportă diferite fonturi (italic, îngroșat, subliniat), precum și imagini, grafice sau diagrame, iar orice document salvat în RTF își păstrează atributele de formatare intacte la transferul între diferite procesoare de text.

Spre deosebire de majoritatea formatelor de procesare de text, codurile RTF pot fi făcute citibile de către om (human-readable). Când se deschide un fișier RTF într-un editor de text, textul alfa-numeric este lizibil și elementele limbajului de formatare (markup language) nu sunt greu de descifrat, lucru ce a constituit o raritate la vremea când acesta a fost lansat. Însă,

fișierele RTF produse de majoritatea programelor, ca MS Word, conțin șiruri atât de mari de numere ale codurilor de control necesare pentru a asigura compatibilitatea cu programele mai vechi, încât cele mai multe din ele au un ordin de mărime mai mare decât textul în sine și sunt foarte greu de citit. Din contră, formatele binare ca și formatul DOC de la MS Word au numai câteva frânturi de text lizibil.

De asemenea, în comparație cu formatul DOC de la Microsoft Word sau cu formatele Open Document (DOT), RTF nu suportă macrocomenzile. De aceea el este recomandat în locul acestor formate pentru că „este protejat și are șanse minime de a fi infectat de virușii informatici. Acest format nu răspândește, deci, viruși când este transmis dintr-un sistem computațional în altul prin poștă electronică. Salvând și trimițând date în format RTF se conferă o siguranță asupra securității fișierului celui care primește documentul.” Totuși extensia RTF nu garantează întotdeauna că un fișier primit este sigur, deoarece Microsoft Word deschide fișierele standard DOC redenumite cu extensia RTF și rulează, astfel, orice macrocomandă conținută în el. De aceea, este necesară examinarea manuală a fișierului într-un simplu editor de text ca Notepad sau folosirea de „file command” în sistemele de tip Unix, pentru a determina dacă un fișier suspect este sau nu de tip RTF.

Fiind recunoscut de mai toate sistemele de operare indiferent de versiunile lor, RTF este un „format util pentru documentele text cu o formatare de bază, cum sunt: diverse manuale de utilizare, rezumate, scrisori și alte documente de informare simple, ce necesită cel puțin o formatare de text (italic, bold sau subliniere), alinierea textului la stânga, dreapta sau central precum și atributele fonturilor și marginilor documentului.”

Aplicația SIL International's Toolbox (<http://www.sil.org/computing/toolbox>), ce dezvoltă și editează dicționare, folosește RTF ca cea mai comună formă a sa de fișier document. Fișierele RTF produse de Toolbox sunt folosite atât în Microsoft Word cât și în celelalte procesoare de text ce recunosc formatul RTF.

Datorită interoperabilității sale, a simplității și a cerințelor mici de procesare a unităților de prelucrare a calculatoarelor (UPC), RTF este un format important pentru instrumentele de citire a cărților electronice (ebook reader) și unele dispozitive, cum este BeBook, lucrează cel mai bine cu acest format.

Programele care deschid fișiere .rtf sunt:

- în Mac OS: Apple TextEditor, Apple Pages, Microsoft Word, Nuance OmniPage ProX sau orice editor de text ce suportă text formatat;
- în Windows: Microsoft WordPad, Microsoft Word, Corel WordPerfect Office X4, Nuance OmniPage Professional 17 sau orice editor de text ce suportă text formatat.

XLS: Microsoft Excel File Format

Microsoft Excel este cel mai utilizat program de calcul tabelar. Un document în format Excel este salvat cu extensia *.xls*.

Excel aparține grupei de programe Microsoft Office și este disponibil pentru Microsoft Windows cât și pentru Mac OS. Actuala versiune disponibilă este pentru *Microsoft Excel 2007* (din 30 noiembrie 2006 pentru firme, respectiv 30 ianuarie 2007 pentru utilizatori) ca și pentru Mac OS *Microsoft Excel 2008* (din ianuarie 2008).

Excel-ul se poate folosi la organizarea datelor, efectuarea de calcule matematice, generarea de rapoarte bazate pe cifre sau pentru a crea diverse grafice.

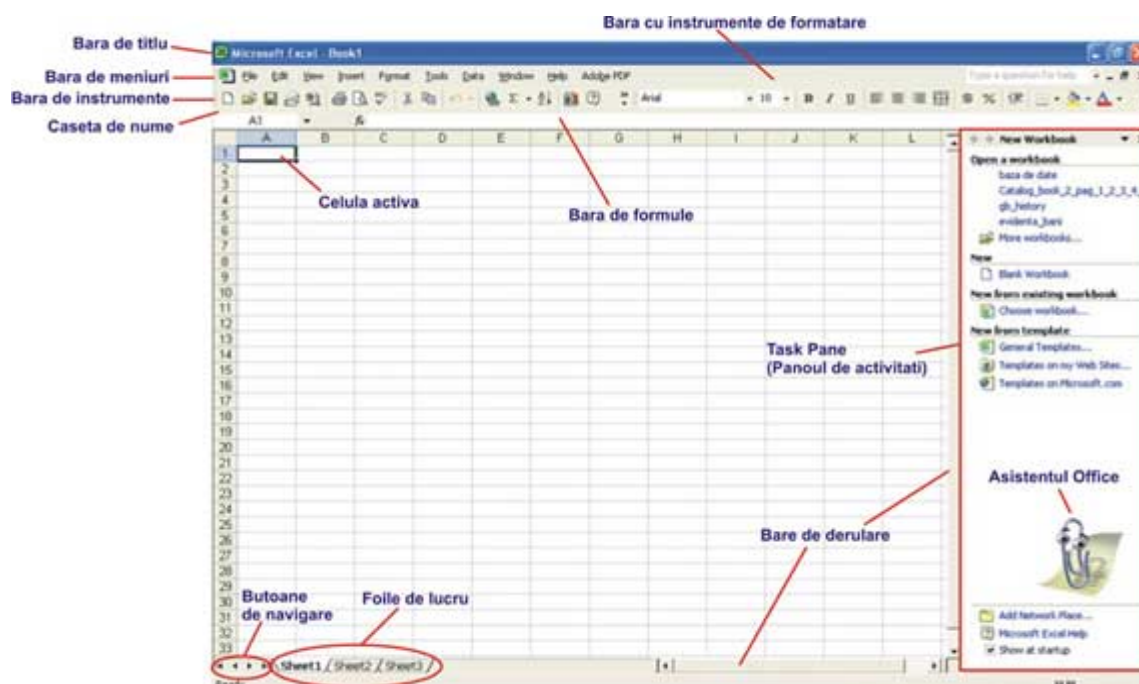


Fig. 2.1. Prezentarea aplicației Microsoft Excel

Tabel 1. Lista de programe care poate deschide fișiere .xls

<i>Nume EXE</i>	<i>Nume produs</i>
Excel.exe	Sistemul Microsoft Office 2007
Excel.exe	Microsoft Office 2000
Excel.exe	Microsoft Excel
Excel.exe	Microsoft Office XP
Excel.exe	Microsoft Office 2003
Moc.exe	Microsoft Open XML Converter
pdfcreator.exe	PDFCreator
scalc.exe	
xlview.exe	Sistemul Microsoft Office 2007

În Microsoft Office 2007, Microsoft prezintă formate noi de fișier pentru Word, Excel și PowerPoint cunoscute ca formate Office Open XML. Aceste noi formate de fișier facilitează integrarea cu surse de date externe și oferă dimensiuni reduse pentru fișiere și operațiuni

îmbunătățite de recuperare de date. În Office Excel 2007, formatul implicit pentru un registru de lucru Excel este formatul de fișier XML Office Excel 2007 (.xlsx). Alte formate XML disponibile sunt formatul XML și cu macrocomenzi activate Office Excel 2007 (.xlsm), formatul de fișier pentru un șablon Excel Office Excel 2007 (.xltx) și formatul de fișier pentru un șablon Excel cu macrocomenzi activate Office Excel 2007 (.xltm).

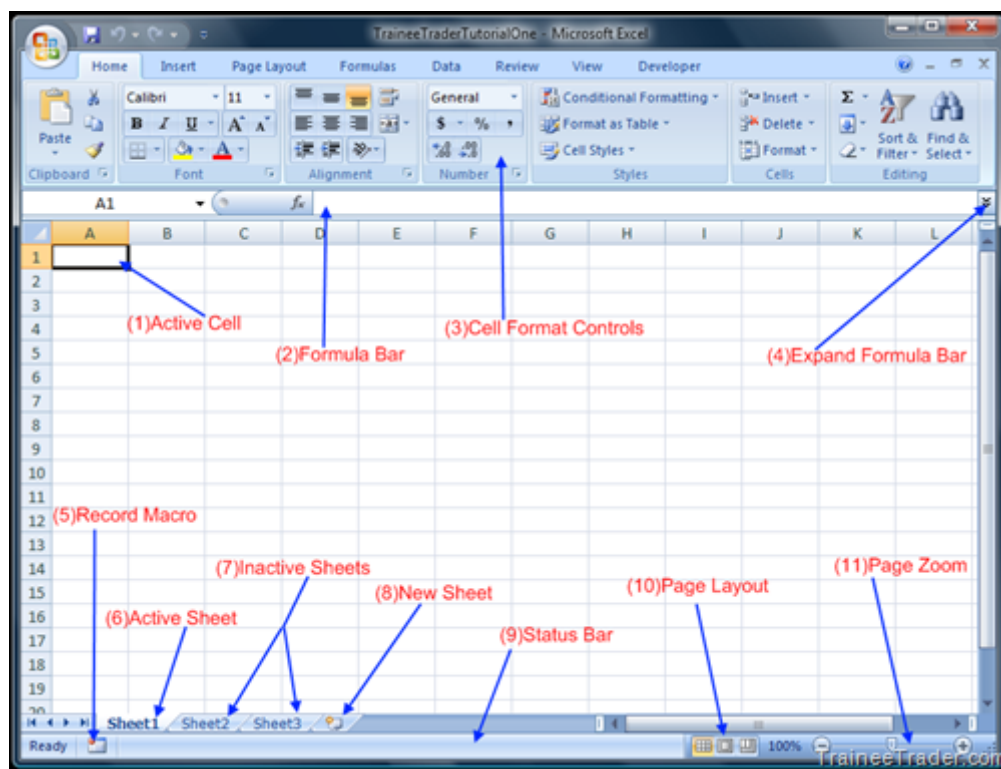


Fig. 2.2. Formatul de fișier pentru un șablon Excel cu macrocomenzi activate Office Excel 2007

Tabel 2. Formate Excel

Format	Extensie	Descriere
Registru de lucru Excel	.xls	Format de fișier binar
Registru de lucru Excel	.xlsx	Formatul de fișier Office Excel. Nu poate stoca cod de macrocomenzi Microsoft Visual Basic pentru aplicații VBA sau foi de macrocomenzi Microsoft Office Excel
Registru de lucru Excel (cod)	.xlsm	Formatul de fișier Office Excel cu macrocomenzi activate. Stocază cod VBA sau foi de macrocomenzi Excel
Registru de lucru binar Excel	.xlsb	Formatul de fișier binar Office Excel (BIFF).
Șablon	.xltx	Formatul de fișier Office Excel pentru un șablon. Nu poate stoca cod de macrocomenzi VBA sau foi de macrocomenzi Office Excel
Șablon (cod)	.xltm	Formatul de fișier Office Excel pentru un

		șablon, cu macrocomenzi activate. Stocază cod de macrocomenzi VBA sau foi de macrocomenzi Excel
Șablon Excel	.xlt	Formatul de fișier binar Excel pentru un șablon Excel.
Foaie de calcul	.xml	Format de fișier foaie de calcul
Program de completare Excel	.xlam	Programul de completare Office Excel cu macrocomenzi activate, un program suplimentar, proiectat pentru a executa un cod suplimentar. Acceptă utilizarea VBA și a foilor de macrocomenzi Excel
Program de completare Excel	.xla	Programul de completare Excel, un program suplimentar proiectat pentru a executa un cod suplimentar; acceptă utilizarea VBA.
Registru de lucru Excel	.xlb	Un format de fișier Excel care salvează numai foi de lucru, foi diagramă și foi de macrocomenzi; există posibilitatea deschiderii unui registru de lucru în acest format de fișier în Office Excel, dar nu și salvarea unui fișier Excel în acest format.

XLSB: Excel binary workbook

XLSB: Registru de lucru binar non-XML

Descrierea extensiei de fișier *.xlsb*:

Este un registru de lucru cu stocare de informații în formă binară și este folosit pentru deschiderea și salvarea documentelor mai rapid și mai eficient, fiind destinat în special pentru documente foarte mari, o cantitate mare de date, cu zeci de mii de rânduri, și/sau de mai multe sute de coloane.

Fișierele cu extensie *.xlsb* pot include tabele de numere, text sau ambele numere și text, formule, conexiuni de date, diagrame și imagini. Conținutul registrului de lucru este de obicei organizat într-o grilă și include date numerice, date structurate și formule.

Cel mai ușor mod de a deschide un fișier *.xlsb* este de a da dublu click pe el, iar computerul decide singur. Dacă fișierul nu se deschide înseamnă că nu este instalată aplicația care poate vizualiza și edita fișiere *.xlsb*.

De asemenea se poate utiliza Notepad sau alt editor de text pentru a deschide fișierul *.xlsb*.

Atunci când un registru de lucru este salvat în formatul *.xlsb*, datele din foaia de calcul sunt memorate în Biff (Binary Interchange Format de fișier). Fiecare parte Biff este salvată ca o

serie de înregistrări Biff și chiar dacă diferitele tipuri de înregistrări Biff conțin informații diferite, fiecare înregistrare are același format de bază. Fiecare înregistrare începe cu un număr de înregistrare urmat de o dimensiune de înregistrare. Aceste două elemente sunt apoi urmate de câmpurile înregistrării. Pentru a converti un fișier *.xlsb* în alt tip de fișier se deschide fișierul *.xlsb* implicit în programul său și se salvează fișierul deschis ca alt tip de fișier sau se folosește un fișier de conversie.

XLSM: Microsoft office excel 2007 document

.XLSM: Registru de lucru care permite macrocomenzi

Descriere fișier

Fișierele cu extensie *.xslm* sunt foi de calcul create cu Microsoft Excel care conțin macrocomenzi activate VBA (Visual Basic for Applications – limbaj de programare derivat din Visual Basic Standard). Doar fișierele care au extensia numelui de fișier terminată în *m* (cum ar fi *.docm* și *.xslm*) pot conține macrocomenzi VBA și controale ActiveX, care sunt stocate într-o secțiune distinctă din fișier. Extensiile de nume de fișier distincte facilitează distingerea fișierelor care conțin macrocomenzi de cele care nu conțin, facilitând identificarea de către software-ul antivirus a fișierelor care conțin cod cu potențial dăunător, prin urmare fișierele cu extensia *.xslm* nu pot fi considerate sigure în cazul în care nu provin de la surse de încredere (pot fi virusate). (O macrocomandă poate fi un mic program sau script care automatizează sarcinile comune. Aceste script-uri rulează în cadrul programelor și sunt create de utilizator.)

XLT: Excel templat

Dacă avem de creat mai multe documente cu o structură asemănătoare și cu informații comune, în loc să le introducem sau copiem de mai multe ori, vom crea un document șablon, pe care-l vom putea folosi la crearea celorlalte documente.

Șablonul se creează ca un document normal dar se salvează cu tipul „Template”; fișierul corespunzător va avea extensia *.xlt*. Pentru a putea fi deschise ușor cu File - New - General, se recomandă salvarea șabloanelor în directorul implicit "Microsoft Office\Templates".

În vederea obținerii documentului creat pe baza unui șablon, se deschide șablonul, se efectuează modificările necesare și se salvează conținutul (ca fișier *.xls*).

XLTM: Excel 2007 template

Un fișier cu extensia *.xltm* este un fișier șablon activat pentru macrocomenzi (macrocomenzile pot să conțină viruși), care poate stoca cod de macrocomenzi VBA (Visual Basic for Applications) sau foi de macrocomenzi Excel 4.0 (*.xlm*).

Este un format nou și se aplică la Word, Excel și Power Point.

ASCII (TXT)

Un fișier **text** este un tip de fișier în care datele sunt stocate ca o secvență de caractere, într-o codificare predefinită (de obicei ASCII, dar mai recent și Unicode). Este unul din cele două tipuri canonice de fișiere, celălalt tip fiind cel de fișier **binar**, în care datele sunt stocate ca o secvență de biți. Diferența dintre cele două tipuri de fișiere este semi - arbitrară, orice fișier text este, până la urmă, un fișier binar, în schimb, nu orice fișier binar este un fișier text. Prin definiție, un fișier text este codificat unitar, conținutul său fiind lizibil și editabil direct, prin intermediul unui editor simplu de text (Notepad, vi, emacs, gedit etc.). Fișierul text este considerat independent de platformă. Într-un fișier text, tipul de date este considerat explicit ca fiind textul neformatat. Organizarea unui astfel de fișier este pe rânduri (care se generează la apăsarea tastei Enter / Return), fiecare rând de text fiind delimitat de unul sau mai multe caractere de sfârșit de rând (caractere EOL). Acestea diferă în funcție de sistemul de operare folosit pentru crearea și editarea fișierului. Astfel, Windows folosește două caractere de control (ASCII 13 urmat de ASCII 10, sau, mai scurt CR+LF) pentru a semnaliza sfârșitul unei linii de text, pe când sistemele tip UNIX (incluzând aici și Linux și Mac OS X) folosesc numai caracterul LF, iar sistemele Mac OS pre-Unix (versiunile 9 sau mai vechi), folosesc doar caracterul CR. Astăzi, majoritatea sistemelor și editoarelor de text sunt perfect capabile să se folosească și să convertească automat fișierele text între diverse arhitecturi. O excepție notabilă în acest sens este editorul Notepad care în cazul în care este pus să deschidă un fișier text generat pe o altă platformă, va afișa caracterul de control LF sau CR după fiecare rând, neîmpărțind textul în rânduri.

Cele mai întâlnite tipuri de codificări sunt:

-**ANSI**: folosit în mod implicit în Notepad, dar este limitat la 256 de caractere; se recomandă pentru caracterele limbii engleze, deoarece o astfel de codificare ocupă mai puțin spațiu și necesită un timp redus de procesare;

-**Unicode**: standard folosit pentru a reprezenta simboluri text din toate limbile, putând recunoaște mai mult de un milion de caractere; include diferite metode;

-**UTF-8**: numărul de octeți pe care este ținut minte un caracter variază de la 1, pentru literele latine, și până la 6 pentru limbile cu accente și diferite simboluri; caracterele nelatine ce erau ținute minte și prin codificarea de tip ANSI sunt procesate mai eficient în acest caz.

Pentru utilizator, un astfel de fișier, vizualizat ca atare, apare ca text neformatat și naranjat în pagină. Fișierele text nu pot stoca decât caractere printabile (litere și cifre) și un număr foarte limitat de caractere de control.

Astfel, în structura fișierului nu se pot stoca elemente de formatare precum culori sau imagini. Fișierele text pot fi folosite pentru activități simple, precum luatul de notițe. Cu toate acestea, un fișier text poate stoca date de intrare pentru un program terț. Astfel de utilizări, datorită portabilității uriașe a formatului (neținând seama de problema sfârșitului de linie) și datorită faptului că conținutul fișierului este ușor lizibil și editabil ca text a dus la apariția multor standarde care codifică informații de diverse feluri ca text lizibil stocat în fișiere text și care procesate cu ajutorul unor alte programe iau altă formă. Exemple în acest sens sunt standardul XML (și aplicațiile sale - HTML - care stochează pagini formatare într-un fișier text prin intermediul unor descriptori textuali sau SVG - care stochează imagini ca text lizibil în fișiere text). Mai mult, codul sursă al oricărui program, în orice limbaj de programare, este salvat pe disc sub formă de fișier text.

Tradițional, fișierele text stochează datele în format ASCII, pe 7 biți. Multe protocoale asociate rețelelor informatice permit transmiterea corectă a informațiilor în acest format, dar nu permit transmisia corectă de fișiere binare (în care codificarea are loc pe 8 biți), ducând, în cazul transmisiei, la coruperea fișierului rezultat (din cauza conversiei ce are loc între cele două sisteme la apariția caracterelor de sfârșit de rând). Astfel, există diverși algoritmi și soluții software pentru codificarea fișierelor binare în format compatibil ASCII pentru transmitere (cel mai comun este Base64). Cu toate acestea, sistemele informatice stochează toate datele într-un mod identic - ca șiruri de biți. Mai mult, unele sisteme de operare nu oferă rutine și proceduri specifice fișierelor text, tratând toate fișierele ca șiruri binare.

Deși acest tip de format nu se mai folosește des, el are avantajele lui:

- dimensiuni reduse datorate lipsei formătării;
- risc redus de a împrăști virusi, conținând doar caractere;
- compatibilitate pe toate platformele;
- nu necesită cumpărarea unui software, deoarece majoritatea editoarelor cu care se deschid sunt încorporate în sistemul de operare.

2.1.2. Formate imagine

Reprezentarea unei imagini pe calculator a constituit un pas firesc, odată cu creșterea puterii de calcul a PC-urilor. O imagine este constituită din pixeli (puncte), denumirea de pixel provenind de la picture elements, el fiind elementul fundamental ce sta la baza unei imagini. Fiecare pixel capătă o anumita culoare dintr-un spectru variabil și este memorat ca un numar, ce va fi înregistrat intern într-o serie de biți. Cu cât numarul de biți alocați este mai mare, cu atât crește posibilitatea de a afișa nuanțe coloristice mai apropiat de realitate.

Pentru o imagine monocrom, fiecare pixel va deține culoarea alb sau negru, deci el va fi memorat ca 0 sau 1. Pentru o imagine în 4 culori, pixelii vor putea lua valori de la 0 la 3 (în format binar 00, 01, 10 și 11 - deci 2 biți), iar pentru 256 de culori (2 la puterea a opta) vor fi necesar câte 8 biți pentru fiecare pixel. Imaginea astfel creată va fi totuși săracă, un numar de culori rezonabil fiind 65536 (16 biti), iar pentru perfecțiune fiind necesare 16 milioane de culori (24 de biti).

Cunoscând faptul ca un byte este format din 8 biți, pentru o imagine pe 24 de biți, un pixel este memorat pe $24/8=3$ bytes. Astfel, pentru o imagine la rezoluția de 1024x768 vor fi necesari $1024 \times 768 \times 3 = 2.359.296$ bytes pe disc. Acesta este spatiul ocupat de o imagine necomprimata, ce poartă de obicei extensia **BMP**, fiind formatul standard, cel mai simplu, dar și cel mai nepractic. În timp, s-a constatat că majoritatea imaginilor nu necesită o fidelitate perfectă și au fost concepuți numeroși algoritmi de compresie, mulți fiind asemănători cu cei folosiți la compresia datelor, dar care duc în multe cazuri la pierderi de calitate. În acest mod, o imagine de 2 MB poate ocupa doar 100 KB în format comprimat, pierderea de calitate fiind în general insesizabila unui ochi amator. În timp, s-a constatat ca majoritatea imaginilor nu necesita o fidelitate perfecta și au fost concepuți numeroși algoritmi de compresie, mulți fiind asemănători cu cei folosiți la compresia datelor, dar care duc în multe cazuri la pierderi de calitate. În acest mod, o imagine de 2 MB poate ocupa doar 100 KB în format comprimat, pierderea de calitate fiind în general ne semnificativă unui ochi amator.

Este important de cunoscut ce format de fotografie trebuie folosit în anumite situații, deoarece unele formate de imagine sunt bune pentru a obține o balanță optimă dintre calitate și mărimea imaginii, în timp ce alte formate sunt perfecte pentru a recupera sau a modifica mult mai ușor caracteristicile unei fotografii (contrast, luminozitate, culoare). Cele mai relevante trei formate în fotografia digitala sunt: Raw, Tiff și JPEG. In continuare vor fi prezentate caracteristicile fiecarui format.

TIFF- Tagged Image File Format

TIFF este un acronim al Tag(ged) Image File Format și totodată unul din cele mai populare și flexibile tip de fișier în format raster. Deși este foarte cunoscut datorită flexibilității și

puterii de compresie, el este considerat în egală măsură destul de complicat și misterios în același timp. Aceasta deoarece formatul TIFF este extensibil și cu multe caracteristici de care un programator ar putea avea nevoie atunci când lucrează cu imagini grafice. Din această cauză el este considerat și ca cel mai confuz format grafic. Tagged Image File Format are mai multe formate interne. Multe programe suportă formatele de bază, însă nu suportă toate formatele existente.

Fișierul TIFF este organizat pe trei mari secțiuni:

1. *Antetul* - În ciuda complexității fișierului TIFF, acesta are un antet foarte simplu. El conține 3 câmpuri și are lungimea de 8 octeți.
2. *Directorul Imaginii* - Este o colecție de informații folosite pentru a descrie imaginea bitmap care urmează. Poate părea puțin cam confuz deoarece Directorul Imaginii poate apărea oriunde în fișier, dar acest lucru nu trebuie să neliniștească deoarece el este pus într-o listă simplu înlănțuită, împreună cu ceilalți Directori, listă spre care pointează Primul ID din Antet. Spre deosebire de structurile întâlnite până acum, structura Directorului Imaginii nu este fixă, ea depinzând de tipul imaginii căruia îi servește ca și antet. În această structură se pot adăuga sau scoate noi câmpuri, exact ca și paginile dintr-un dosar.
3. *Zona de Date* în care sunt stocate imaginile bitmap - Fișierele TIFF contin numai imagini de tip bitmap, deși cu ajutorul câmpurilor nu este exclusă nici stocarea unor imagini vectoriale, sau a unor texte (vezi și formatul GIF).

Din aceste trei secțiuni este nevoie doar de primele două, existând posibilitatea ca Zona de Date să lipsească. TIFF are reputația de a fi un format complicat datorită faptului că în interiorul fișierului numai antetul poate fi găsit la o locație fixă, restul secțiunilor variind de la fișier la fișier. Fiecare Director al Imaginii împreună cu bitmap-ul asociat formează un subfișier. Nu există o limită a numărului de subfișiere pe care le poate conține un fișier TIFF.

Formatul TIFF suportă cei mai mulți algoritmi de compresie. Versiunea 4.0, are incluși numai algoritmi Run Length Encoding, versiunea 5.0 adaugă la acestea și algoritmul LZW, întâlnit și la formatul GIF, iar versiunea 6.0 adoptă și algoritmul JPEG, algoritm care se descurcă de minune cu împărțirea fișierului în cărămizi (o bucată dreptunghiulară din imaginea bitmap, care a fost introdusă datorită imaginilor TIFF foarte mari).

Totodată TIFF este un format comun utilizat pentru diverse aplicații pentru imagini, inclusiv cele pentru scanare și fax. Microsoft Office Document Imaging utilizează formatul

TIFF, folosind posibilitatea acestuia de a conține text acceptat de recunoașterea optică a caracterelor (OCR) (OCR: traduce imaginile de text, cum sunt documentele scanate, în caractere de text efective. Este cunoscut și ca recunoaștere de text.). Când se scanează documente noi, acestea sunt salvate în format TIFF (cu extensia tif) și orice text OCR este stocat în fișierul TIFF împreună cu imaginea.

Se pot deschide și edita fișiere TIFF create cu Office Document Imaging în multe alte aplicații grafice. În acest caz, orice text OCR conținut în fișier, va fi pierdut. Dacă se dorește să se acceseze din nou textul din fișierul TIFF, în Office Document Imaging, va trebui să se reexecute OCR.

Detalii tehnice:

Office Document Imaging creează fișiere în următoarele formate:

Monocrom Un bit per pixel, compresie G4

Tonuri de gri 8 biți per pixel, compresie JPEG

Color 24 biți RGB, compresie JPEG

Office Document Imaging acceptă:

- Toate tipurile de compresii listate în specificația TIFF 6.0.
- Diferite tipuri de compresii pentru fiecare pagină a unui document multipagină.
- Imagini TIFF cu adâncimea culorii de 1-bit, 4-biți, 8-biți sau 24-biți (atât paletă cât și non-paletă).
- Spațiere culoare RGB și CMYK.
- Imagini rearanjate și redimensionate.

Office Document Imaging nu acceptă:

- Spațiere culoare YCbCr, cu excepția imaginilor JPEG.
- Spațiere culoare CIE Lab.
- Imagini cu mai mult de cinci eșantioane per pixel sau cu o dimensiune a eșantionului mai mare de 32 biți.
- Imagini în format Planar.

TIFF este un standard care se folosește în mod special în industria tipografică. Fișierele TIFF sunt semnificativ mai mari decât cele JPEG și pot fi și ele necomprimate sau comprimate folosind compresia lossless. Spre deosebire de JPEG, fișierele TIFF pot fi de 16-biți/channel sau 8-biți/channel și pot fi folosite layer pentru a stoca mai multe imagini într-un singur fișier TIFF.

Fisierele TIFF sunt o opțiune excelentă de backup și de a fi folosite ca fișiere intermediare, care pot fi editate mai târziu, din moment ce nu sunt folosite artefacte de compresie. Multe camere de fotografiat au opțiunea de a folosi fișiere TIFF, dar acestea pot consuma un spațiu destul de mare și destul de important de pe cardul de memorie a aparatului. În cazul în care camera suportă formatul RAW este indicat folosirea acestuia deoarece este o alternativă superioară față de TIFF.

Avantajele formatului:

- Format universal, recunoscut de majoritatea programelor.
- Calitate foarte bună a fotografiilor.
- Compresie (opțională) prin protocol LZW.
- Format ideal pentru stocarea imaginilor în curs de prelucrare sau prelucrate.
- Codificare până la 64 biti.
- Permite gestionarea transparenței (strat alfa).
- Conversie în orice alt format fără pierdere de calitate.

Dezavantajele formatului:

- Fișierele de imagine au dimensiuni foarte mari.
 - Nu este adecvat pentru utilizarea în camerele foto digitale.
 - Nepotrivit pentru utilizare pe Internet.

GIF- Graphics Interchange Format

Formatul Gif a fost creat pentru elemente grafice de dimensiuni mici. Modul de lucru se bazează pe reducerea paletelor de culori a graficii la maximum 256 de culori. Adesea, ceea ce se percepe ca fiind o culoare sunt de fapt mai multe: de exemplu, galbenul este compus din mai multe nuanțe sau, în cazul culorilor digitale, din pixeli de culori individuale.

Tehnologia GIF reduce numărul de culori folosind *intercalarea* pentru cele care nu se găsesc în paleta GIF standard de 256 de culori.

GIF sau *Graphics Interchange Format* (Format de schimb de grafică), este formatul de fișier folosit de cele mai multe persoane pentru a face schimb de grafică. Devenit popular la origine pe CompuServe, GIF s-a răspândit către alte servicii on-line și apoi către Internet și către Web. Orice browser ce suportă grafică, suportă GIF.

Pentru imaginile cu o mulțime de blocuri de culoare uniformă, dimensiunile fișierelor GIF tind să fie mici. Astfel, se preferă GIF-urile pentru bannere (steaguri) sau pentru imaginile cu zone întinse de culoare uniformă, cum ar fi graficele cu bare sau pictogramele. Deci desenele simple pe care le creează cei mai mulți dintre noi funcționează cel mai bine cu GIF. Grafica artistică densă și fotografiile funcționează mai bine cu JPEG.

Fișierele GIF dau, de asemenea, anumite opțiuni de afișare a paginii Web care nu se obțin întotdeauna cu fișierele JPEG. Se pot face culorile din imaginile GIF transparente față de orice se află în fundalul imaginii și se pot salva imaginile GIF în format *intercalated* (*întrețesut*). Imaginile salvate în acest mod și apoi descărcate (download-ate) de către un browser, apar mai întâi într-o rezoluție foarte scăzută și apoi într-o rezoluție progresiv mai clară, până apare întreaga imagine. Această facilitate face imaginile GIF preferabile pentru afișarea rapidă a unei grafici brute care se îmbunătățește pe măsura trecerii timpului și pentru crearea efectelor speciale extravagante.

Este desigur unul dintre cele mai cunoscute formate pentru fișiere de imagini: GIF (*Graphic Interchange Format*) a fost dezvoltat de CompuServe în 1987 și este foarte des utilizată pe WWW grație portabilității sale. Formatul utilizează o paletă până în 256 de culori din cei 24 biți de culoare din spațiul RGB. Extensia acestor fișiere este .GIF.

Suportă animații și alocă câte o paletă de 256 de culori pentru fiecare frame. Poate fi folosit pentru reproducerea de filmulețe cu rezoluții mici. Limitările de culoare fac ca formatul Gif să fie nepotrivit pentru reproducerea pozelor sau altor imagini cu culori continui, dar este foarte potrivit pentru imagini simple ca exemplu grafice sau logo-uri cu suprafețe de culori unanime.

Imaginile gif utilizează un algoritm de compresie LZW - denumit după cei care l-au creat, Lempel, Zif și Welch în 1985- această tehnică permite reducerea dimensiunii imaginii fără degradarea vizuală a imaginii. La imagini color este posibilă o comprimare de 3:1, uneori se poate obține chiar și o rată de comprimare de 5:1. Codificarea LZW divizează informațiile de imagine. Fiecărei succesiuni îi este atribuit un index, care este salvat într-un tabel. În consecință, informația propriu-zisă nu este formată din puncte individuale de imagine, ci din indexurile din tabel.

La prima vedere, comprimarea LZW lucrează fără pierderi în cazul fișierelor GIF, deoarece fișierul GIF decomprimat este identic cu cel dinainte de comprimare. Problema este însă următoarea: fișierele GIF salvează numai 8 biți pe punct de imagine, fapt ce micșorează

adâncimea de culoare la 256 de culori. Dar, fiind vorba despre plăci grafice care pot reprezenta până la 16 milioane de culori, această adâncime este prea mică.

Cu ceva timp în urmă, GIF a fost cel mai important format de fișier în Usenet pentru transferul de imagini color prin intermediul liniei telefonice. Însă, pentru ca Unisys a protejat compresia LZW prin patent, au existat tot mereu probleme de licență la aplicațiile comerciale. Acest fapt și metodele moderne de comprimare care au apărut între timp au dus la înlocuirea lui GIF cu JPEG.

Mai vechi decât JPEG, pentru unele tipuri de imagini, GIF este superior în privința calității imaginii, a ratei de comprimare sau a ambelor, dar are un defect: o paletă redusă de culori, deoarece, prin convenție, GIF are maximum 256 culori (codificare pe 8 biți). De exemplu, arii largi cu exact aceeași culoare, sunt comprimate foarte eficient de algoritmul GIF. Dar, pentru o imagine fotografică, transpusă digital prin scanare sau înregistrată cu aparate fotografice digitale, adâncimea de culoare este de cel puțin 24 biți (16 milioane de culori), ceea ce determină o alterare severă a pozei.

Imaginile optime pentru comprimare prin metoda GIF sunt cele cu linii și cu suprafețe uniform colorate, cum sunt cele desenate cu programe de grafică (de exemplu Microsoft Paint). Cu cât imaginea conține mai multe nuanțe, cu atât JPEG se descurcă mai bine. GIF comprimă bine marginile netede, cum ar fi chenarele fotografiilor sau conturul literelor aplicate peste o imagine, pe care JPEG le redă mai difuz, datorită rotunjirilor inerente în calculele matematice pe care le efectuează algoritmul de compresie.

În plus, formatul GIF permite animații simple. Formatul GIF este foarte folosit pentru realizarea site-urilor Internet (butoane, bare de defilare, logo-uri, etc).

Avantaje

- Format universal, fără copyright, foarte folosit pentru Internet, recunoscut de orice navigator, fără să aibă nevoie de plugin-uri;
- Comprimă imaginea fără pierderi de calitate, până la dimensiuni foarte mici;
- Posibilitate de animare;
- Poate gestiona transparența.

Dezavantaje

- Codificare doar pe 8 biți;
- Reproducere de calitate redusă a imaginilor fotografice.

JPEG- Joint Photographic Experts Group

Formatul de fișier JPEG (*Joint Photographic Experts Group*) a fost dezvoltat de C-Cube Microsystems în 1992 și aprobat în 1994 ca standardul ISO 10918-1. În ciuda adâncimii mari de culoare (16 milioane de culori) imaginile scanate cu JPEG pot fi stocate pe harddisk, ocupând puțin spațiu. Acest lucru este posibil datorită unei comprimări a datelor bine pusă la punct, care micșorează dimensiunile fișierului la minim. Fișierele JPEG pot fi recunoscute după extensia: .jpeg, .jpg, .jpe, .jfif, .jfi, .jif (containers).

Formatul JPEG folosește pentru comprimare transformarea cosinus discretă (DCT). În spatele acestui sistem se ascunde o transformată Fourier, care schimbă dispersia pixelilor într-o dispersie de frecvență de amplitudine. Suprafețelor mari și uniforme de imagine li se atribuie participării de frecvență mici, în schimb detaliile fine primesc participării de frecvență mai mari. Comprimarea apare ca urmare a faptului că părți de imagine cu participării de frecvență mai mari capătă o pondere mai mică, iar amplitudinile lor vor fi egale cu zero.

La comprimare se pierd unele date de imagine, însă algoritmul de compresie utilizat verifică în permanență că la decomprimare imaginea să nu sufere o pierdere a calității. La comprimare, fiecare imagine este divizată în blocuri de 8x8 puncte de imagine. Pentru început, fiecare bloc trebuie să suporte o transformare cosinus. Scopul acesteia este de a cuprinde modificarea de culoare de la un bloc la altul. În etapa a doua, algoritmul verifică dacă diferențele de culoare dintre blocuri sunt vizibile cu ochiul liber. Diferențele dintre blocuri vecine care nu sunt observabile cu ochiul liber sunt ignorate.

În ultima etapă vor fi comprimate celelalte blocuri în mod obișnuit și vor fi stocate pe harddisk într-un fișier JPEG. În funcție de conținutul imaginii, această metodă economisește între 50 și 70 de procente din spațiu de stocare, fără să fie observată o pierdere a calității.

O consecință a comprimării puternice cu JPEG este faptul că aceste fișiere pot fi transferate prin modem prin intermediul liniei telefonice. Astfel, majoritatea imaginilor din Internet sunt codificate JPEG. JPEG aduce cu sine, în schimb, și o serie întreaga de dezavantaje. Dacă în cazul fotografiilor comprimarea decurge fără probleme, în schimb la desenele tehnice, de exemplu, apar probleme serioase.

O proprietate foarte utilă a algoritmului JPEG este capacitatea acestuia de a avea un grad variabil de comprimare, ales de utilizator. Aceasta înseamnă că dacă se dorește obținerea unui fișier de imagine cât mai mic, se poate alege o rată mare de comprimare, în dauna calității; invers, pentru a menține calitatea la o cotă ridicată, se alege un grad redus de comprimare. La prima comprimare a imaginilor, chiar și pentru un grad mediu, pierderea de calitate este minoră. Așa se face că toate aparatele fotografice digitale, pentru a realiza economie de spațiu pe cartela de memorie, prezintă posibilitatea de a comprima imaginea descărcată de pe CCD în imagine comprimată JPEG în diverse rapoarte (1:5 - 1:20), fără o alterare semnificativă. Dacă însă fișierul JPEG este deschis, prelucrat, salvat, apoi din nou deschis, prelucrat ... de mai multe ori, pierderile de calitate devin tot mai evidente. Cât de mult se poate comprima o imagine? Dacă se pleacă de la un fișier imagine de tip TIFF de 1 MB (să presupunem), putem obține un fișier jpg de 100 K (rata de compresie 1:10), fără alterări importante ale calității vizibile a imaginii. O comprimare la o rată de 1:40 - 1:50 aduce fișierul-imagine la 20 - 25 K cu pierderi perceptibile dar moderate. Este posibilă comprimarea de până la 1:1000, imaginea obținută având doar 1 K și păstrează mai multe detalii decât oferă multe programe de manipulare a imaginilor în modul "thumbnail".

Avantajele formatului

-Format universal, fără copyright, recunoscut de orice program de vizionare și prelucrare a imaginilor;

-Foarte bun raport calitate+dimensiuni;

-Cel mai folosit format pentru arhivarea fotografiilor;

-Utilizat aproape exclusiv pentru afișarea fotografiilor pe internet;

-Codificare pe 24 biți (până la 16 milioane de culori).

Dezavantajele formatului

-Format cu pierderi, datorită algoritmului de compresie;

-Fișierele prea mari pentru grafice;

-Nu gestionează straturi Alfa (transparență).

Fișiere gif întrețesute și fișiere jpg progresive

Randarea progresivă reprezintă apariția gradată a graficii, fiind o metodă utilă pentru a păstra atenția pe durata încărcării. Imaginea apare mai întâi “încețoșată” și apoi, în mod progresiv, devine mai clară.

Există două metode de randare progresivă în Web. Cea mai populară și mai eficientă este folosirea fișierelor GIF întrețesute. Procedul se poate aplica oricărui fișier GIF și, de obicei, în programele grafice sau instrumente plug-in există opțiuni care să ajute.

Fișierele JPG *progresive* reprezintă răspunsul JPEG la întrețesere. Nu se poate face întrețeserea unui fișier JPG, dar există programe care creează ceea ce este cunoscut sub numele de fișier JPG progresiv. Exemplu: Photoshop 4.0 are opțiunea de randare progresivă pentru JPG.

BMP- bitmap

De când există PC-urile, informațiile de imagine se salvează în așa-numitele fișiere bitmap. Sub sistemul de operare Windows, fișierele de format BMP se folosesc pentru a stoca imagini, icon-uri, cursoare, pointeri, etc.

Acestea pot fi recunoscute datorită extensiei *.BMP. Faptul că acest format, care ocupă mult spațiu de stocare, a supraviețuit atâta timp este de înțeles: în fișierele BMP sunt stocate informațiile aproape așa cum Windows reprezintă intern imaginile. Aici, informația de culoare a fiecărui punct de imagine este stocată în 1, 4, 8, 16 sau 24 biți. Cu cât sunt rezervați mai mulți biți de culoare pentru fiecare punct de imagine, cu atât este mai mare adâncimea culorii și implicit calitatea imaginii. Acest format prezintă însă și un dezavantaj vizavi de celelalte: în măsura în care este mai mare adâncimea de culoare, crește și spațiul ocupat de fișier pe harddisk. Astfel se adună rapid câțiva MB, chiar și în cazul unor imagini color mici, la o rezoluție de 640x480 pixeli. Fiecare fișier BMP este compus din mai multe părți:

- un header;
- un tabel de culoare (opțional);
- datele pentru punctele de imagine.

În header se află, de exemplu, informații despre câte puncte are înălțimea și lățimea imaginii și unde sunt stocate în fișier informațiile de imagine. În tabelul de culoare este stabilit modul în care sunt reprezentate culorile pe paleta de culori a plăcii grafice. Acest lucru este important mai ales atunci când imaginea este salvată cu 16 milioane de culori, însă placa grafică

poate reprezenta numai 256 de culori. La o adâncime de culoare de 16 și 24 de biți nu există un tabel de culoare, deoarece în acest caz informația de culoare este stocată direct în bitmap.

Deci informațiile de imagine sunt stocate într-un fișier BMP în moduri diferite. Utilizatorul nu trebuie să-și bată capul cu aceste amănunte, deoarece filtrele de import ale tuturor programelor de prelucrare de imagini afișează corect informațiile de imagine.

Fișierele BMP pot fi transferate fără probleme între calculatoare. Deoarece datele nu sunt comprimate, la reprezentarea pe monitor sau la imprimare nu are de suferit calitatea imaginii. Dar cei care vor dori să colecționeze pe harddisk mai multe imagini, este bine să evite fișierele BMP. În caz contrar se vor trezi rapid cu un harddisk plin.

Pixelii dintr-o imagine de tip bitmap sunt comprimați folosind o combinație a trei moduri de lucru:

- *Modul RLE* (run length encoded), transformă extensia din BMP în RLE, este de 2 octeți și reprezintă de la 1 la 255 de pixeli, toți de aceeași culoare. În secvența 04 07, de exemplu, urmează patru pixeli având culoarea 07;
- *Modul escape*, primul octet este 0 iar următorii octeți semnifică una dintre următoarele trei posibilități: 0 reprezintă sfârșitul unei linii de imagine, 1 reprezintă sfârșitul fișierului bitmap, 2 indică o comandă delta. O comandă delta mai este urmată de încă doi octeți care reprezintă deplasarea pe orizontală și pe verticală privind locul unde va apare următorul pixel, în raport cu pixelul curent;
- *Absolute mode*. Prin acest mod primul octet are valoarea zero iar următorul octet are o valoare de 3 sau mai mare reprezentând numărul de pixeli necomprizați care urmează. Un exemplu de folosire a modului absolut este 00 03 09 08 06 care interpretat în mod adecvat ne spune că urmează trei pixeli având valorile 09, 08 și 06.

PNG- Portable Network Graphics

Înca din 1994 a apărut ideea de format grafic nou, fără pierderi de calitate, care să suporte 24 de biți de culoare și mai ales, să fie gratuit. Astfel, în octombrie 1996 s-a născut **PNG** (Portable Network Graphics), un format excepțional pentru acea vreme, dar care s-a impus destul de greu. Avantajul sau cel mai important erau pierderile de calitate nule, iar ca dezavantaj putem enunța viteza destul de scăzută a algoritmului de compresie. PNG este folosit în special pentru compresia imaginilor schematice, dar și a celor de tip fotografie unde nu se acceptă pierderi de calitate. Dacă pentru mult timp au fost preferate alte formate, cum ar fi GIF sau JPG, din simplul motiv al compatibilității și - de ce nu - al comodității, în prezent, orice aplicație din domeniu

poate lucra cu acest format. Odată cu introducerea suportului în browsere (fie direct, fie prin diverse plugin-uri), succesul i-a fost asigurat, el fiind în prezent cel mai performant format grafic de pe piață din sfera Internet-ului, însă totuși nu este atât de răspândit ca GIF sau JPG. „Norocul” său a fost faptul că prin definiție, el este destinat web-ului, chiar dezvoltatorii lui numindu-se The World Wide Web Consortium (W3C). PNG este unul din puținele cazuri din istorie în care un format mai performant dar initial necunoscut, a putut cuceri piata. Să ne amintim de muzica și de formatul MP3, care este în prezent este destul de slab ca performanțe, dar extrem de popular, alternativele chiar dacă sunt mai bune din toate punctele de vedere, sunt ignorate de piață,.

În încercarea de a corecta deficiențele fișierelor de tip GIF, a aparut formatul PNG (Portable Network Graphic) care prezintă codificare pe 8 biți (PNG-8) sau pe 24 biți (PNG-24), inclusiv gestionarea transparenței pe 256 nivele. Formatul PNG este recunoscut de practic toate programele moderne de prelucrare și navigare pe Internet. Pentru aceeași imagine, fișierele PNG sunt mai mici decât cele în varianta GIF. Formatul PNG-24 suportă și imagini cu tonuri continue, dar dimensiunile fișierelor sunt mai mari decât cele JPEG.

Avantaje:

- Format universal, folosit pentru Internet, recunoscut de orice navigator modern, fără să aibă nevoie de plugin-uri.
- Comprimă imaginea fără pierderi de calitate, până la dimensiuni foarte mici. Poate gestiona transparența mai bine decat formatul GIF.

Dezavantaje:

- Incompatibil cu versiunile vechi de navigare Internet.
- Fotografiiile sunt comprimate mai puțin decât formatul JPEG.
- Fara posibilitate de animare.

DIB- device independent bitmap

DIB nu este un format de fișier imagine, ci este formatul în care pot fi păstrate imaginile în memorie de aplicațiile Windows. De obicei programele care doresc afișarea diferitelor formate de fișiere imagine recurg la acest format ca și format intermediar. DIB reprezintă o hartă de pixeli, independentă de dispozitiv.

Tabloul de pixeli poate fi memorat necomprimat sau folosind o metodă de compresie de tip RLE (Run-Length Encoded) pe 4 sau pe 8 biți. Metoda de stocare a valorii pixelilor este specificată de valoarea câmpului biCompression din headerul formatului.

În formulare și rapoarte există posibilitatea să se afișeze imagini, precum și ilustrații, sigle sau fotografii. Pentru aceasta, imaginile trebuie mai întâi stocate. Access furnizează câteva metode de stocare a imaginilor.

O posibilitate este următoarea: încorporarea imaginilor într-un câmp de tip Obiect OLE dintr-un tabel al bazei de date.

OLE este tehnologia utilizată pentru a partaja fișiere între diversele programe Office. De exemplu, când se inserează o foaie de calcul Excel într-un document Word sau se inserează un diapozitiv Microsoft PowerPoint într-un desen Microsoft Visio, se utilizează OLE. Un câmp Obiect OLE se utilizează când este necesară stocarea imaginilor (sau legături la ele) și a altor fișiere din alte programe Office direct în baza de date.

Această metodă este cel mai simplu de implementat deoarece se utilizează ecranele și instrumentele furnizate de Access. De asemenea, imaginile devin parte a bazei de date și se deplasează împreună cu ea. Nu va fi nevoie niciodată să se actualizeze legăturile la fișierele imagine.

Încorporarea imaginilor, însă, poate duce rapid la creșterea dimensiunii bazei de date și la o execuție lentă. Aceasta se întâmplă în special când se stochează fișiere GIF și JPEG, deoarece OLE creează fișiere bitmap suplimentare care conțin informații despre afișare pentru fiecare fișier imagine, iar aceste fișiere suplimentare pot fi mai mari decât imaginile inițiale. În plus, această metodă acceptă numai formatele de fișier grafic Windows Bitmap (.bmp) și Device Independent Bitmap (.dib). Dacă se dorește să se afișeze alte tipuri obișnuite de fișiere imagine, cum ar fi imaginile GIF și JPEG, trebuie să se instaleze software suplimentar.

GPX- GPX exchange format

GPX este un format de prezentare a datelor GPS pentru aplicații software derivat din XML (Extensible Markup Language). Poate fi folosit pentru a descrie coordonate și rute. Formatul nu este supus vreunei licențe și poate fi folosit fără a plăti vreo taxă.

Marcajele pe care le folosește stochează locații, altitudini sau date temporale și astfel GPX poate fi folosit pentru schimbul de date între echipamentele GPS și aplicațiile software.

Acest aplicații pot fi folosite împreună cu programe de imagistică prin satelit, cum ar fi Google Earth.

În GPX, un grup de puncte fără vreo legătură secvențială între ele este tratat ca un grup de coordonate. Un grup ordonat de puncte poate fi tratat ca o rută sau ca un traseu. Un traseu reprezintă înregistrarea locurilor pe unde a trecut o persoană, pe când rutele reprezintă destinații posibile în viitor. Așadar, în cazul traseelor pot exista marcaje temporale pentru fiecare punct din traseu (pentru că cineva a notat când a trecut pe acolo), dar în cazul rutelor este puțin probabil ca acestea să fie furnizate. Datele minimale care trebuie incluse într-un document GPX sunt latitudinea și longitudinea pentru o coordonată, restul variabilelor fiind opționale.

EPS- Encapsulated postscript

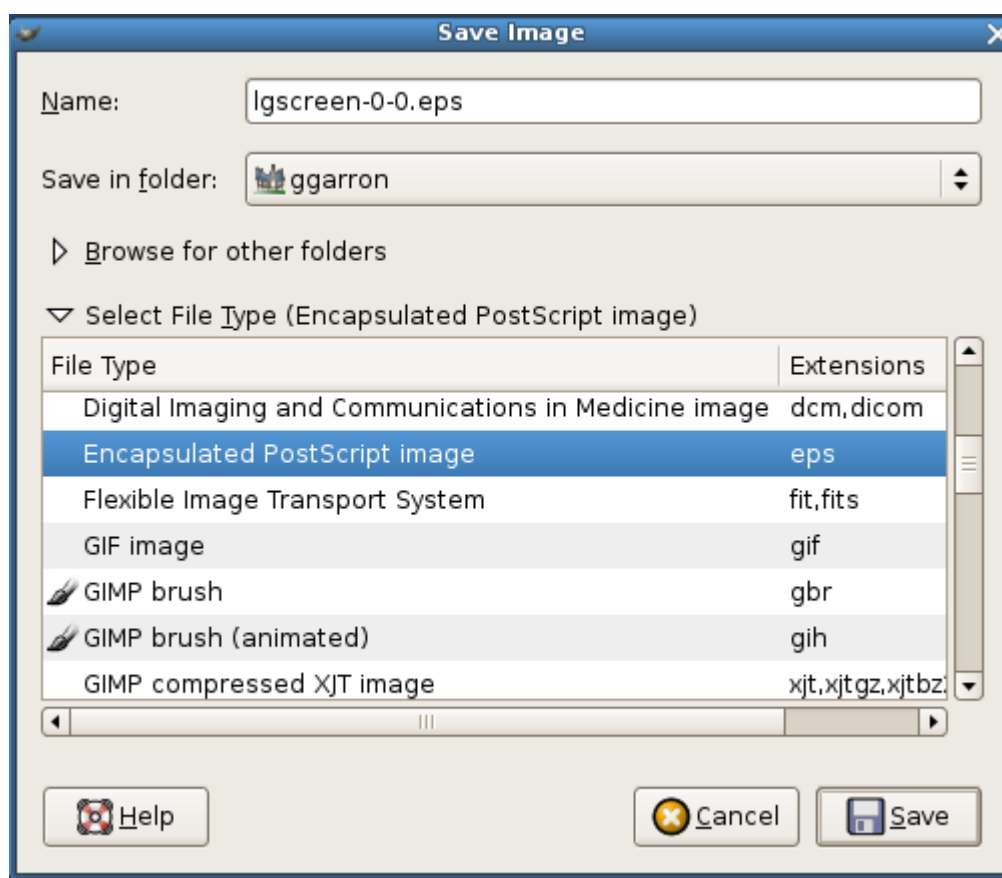


Fig. 2.3. Encapsulated postscript

Formatul de fișier Encapsulated PostScript (*EPS*) poate conține atât grafică vectorială cât și bitmap și este acceptat teoretic de toate programele pentru grafică, ilustrații și machetare a paginilor. Formatul *EPS* este utilizat pentru transferul ilustrațiilor PostScript între aplicații. Când se deschide un fișier *EPS* care conține grafică vectorială, Photoshop rasterizează imaginea, convertind grafica vectorială în pixeli. (Rasterizarea = convertirea imaginilor și textului unui document în imagini raster - un model electronic constând în pixeli - pentru a fi tiparite sau

printate. Aplicația RIP transformă informațiile vectoriale în instrucțiuni executabile de echipamentele de tipărire.)

Formatul *EPS* acceptă modurile de culoare Lab, CMYK, RGB, Culoare indexată, Bicromie, tonuri de gri și Bitmap, și nu acceptă canalele alfa. *EPS* nu acceptă traseele de tăiere. Formatul Desktop Color Separations (DCS) - o versiune a formatului *EPS* standard care permite salvarea separațiilor de culoare ale imaginilor CMYK. Pentru a tipări fișiere *EPS*, trebuie utilizată o imprimantă PostScript.

Photoshop utilizează formatele EPS TIFF și EPS PICT pentru a permite deschiderea de imagini salvate în formate de fișiere care creează previzualizări, dar care nu sunt acceptate de Photoshop (cum ar fi QuarkXPress).

Encapsulated PostScript poate avea extensiile: *.eps* , *.epsf* (grafice și imagini porționate ale documentului), *.epsi* (*bitmap*).

CDR: Corel draw format

CorelDRAW este un program de grafică vectorială, creat pentru a rula pe platforme Microsoft Windows (cel mai important pachet de grafică sub Windows). El permite crearea unor ilustrații mai deosebite, publicitate la nivel profesional, rapoarte anuale. Este un program de desenare orientat pe obiecte. Din multe puncte de vedere este asemănător cu programele de proiectare asistată pe calculator (AutoCad). Acest program se aseamănă cu programele de pictură care folosesc grafica realizată prin puncte și dau iluzia picturii pe pânză. Programele bazate pe grafică vectorială, cum este CorelDRAW sunt de preferat celor ce realizează grafica prin puncte. Avantajele graficii vectoriale sunt precizia și flexibilitatea.

Tabel 3. Fișiere recunoscute de CorelDRAW

Adobe Illustrator	(AI)
Adobe Photoshop	(PSD)
Adobe Type 1 Font	(PFB)
ANSI Text	(TXT)
AutoCAD Database	(DWG)
AutoCAD Interchange Format	(DXF)
CALs Compressed Bitmap	(CAL)
Computer Graphics Metafile	(CGM)
Corel ArtShow 5	(CPX)
Corel DESIGNER	(DSF or DES)
Corel DESIGNER	(CDT)
Corel Paint Shop Pro	(PSP)
Corel Painter	(RIF)
Corel PHOTO-PAINT	(CPT)
Corel Presentation Exchange	(CMX)
Corel Presentations	(SHW)
Corel R.A.V.E.	(CLK)
CorelDRAW	(CDR)

CorelDRAW Compressed	(CDX)
Encapsulated PostScript	(EPS)
FPX	
Frame Vector Metafile	(FMV)
GEM File	(GEM)
GEM Paint	(IMG)
GIF	
GIMP	(XCF)
Hewlett-Packard Plotter	(PLT)
HTML	
JPEG	(JPG)
JPEG 2000	(JP2)
Lotus PIC	(PIC)
MacPaint Bitmap	(MAC)
Macromedia Flash	(SWF)
Macromedia FreeHand	(FH)
MET Metafile	(MET)
Micrografx 2.x, 3.x	(DRW)
Micrografx Picture Publisher 4	(PP4)
Micrografx Picture Publisher 5	(PP5)
Microsoft PowerPoint	(PPT)
Microsoft Word Document	(DOC sau RTF)
NAP Metafile	(NAP)
OS/2 Bitmap	(BMP)
Pattern File	(PAT)
Macintosh PICT	(PCT)
Picture Publisher File	(PPF)
Portable Document Format	(PDF)
Portable Network Graphics	(PNG)
PostScript	(PS sau PRN)
PostScript Interpreted	(PRN)
PostScript Interpreted	(PS)
Rich Text Format	(RTF)
Scalable Vector Graphics	(SVG)
SCITEX CT Bitmap	(SCT)
TIFF Bitmap	(TIF)
TrueType Font	(TTF)
Visio	(VSD)
Windows Bitmap	(BMP)
Windows Metafile Format	(WMF)
WordPerfect Document	(WPD)
WordPerfect Graphic	(WPG)
XPixmap Image	(XPM)

CMX: Vector image file format

CMX: Corel Presentation Exchange Image

Cmx este un format vectorial grafic, un fișier *.cmx* conținând o descriere a modului în care se trasează imaginea și nu doar o descriere a modului de afișare. Apoi imaginea poate fi afișată la orice dimensiune, fiind interpretată, scalată și afișată de către browser-ul plug in.

Fiecare obiect într-o imagine vectorială este stocat ca un element distinct, cu informații despre poziția sa din imagine, dimensiune, culoare etc. O imagine în format vectorial este de rezoluție-independentă. Ea poate fi redimensionată, fără a pierde nici un detaliu, pentru că este stocată ca un set de instrucțiuni, nu o colecție de pixeli.

Acest format nu este potrivit pentru fotografiile obișnuite, din “viața reală”, el funcționând foarte bine pentru planuri, hărți, miniaturi (clipart) și grafice de afaceri. Imaginile *.cmx* pot fi vizualizate numai de Netscape Navigator TM 2.0 pentru Windows 95 și Windows NT.

RAW: Raw image format

RAW (Raw image format) este un format de imagine în care sunt stocate datele brute, neprocesate, capturate de senzorii unei camere foto digitale sau ale unui scanner.

Atunci când au apărut scannerele și camerele foto digitale, senzorii nu aveau calitatea necesară pentru a captura prea multe nuanțe de culoare, așa că spațiul de culoare sRGB era suficient pentru stocarea imaginii capturate. Pe măsură ce calitatea senzorilor s-a îmbunătățit, dispozitivele de captură au ajuns să poată capta o plajă mai mare de culori și nivele de luminozitate decât cea care poate fi stocată într-un fișier JPEG, de exemplu. Ca atare, cineva trebuia să decidă în ce fel să fie expusă imaginea: oricât de bun ar fi fost senzorul, imaginea finală din spațiul sRGB sau Adobe RGB trebuia în mod necesar să reprezinte o gamă dinamică mai mică decât cea capturată de senzori. Cele mai multe dispozitive moderne de captură de calitate rezonabilă, atât scannere cât și camere de fotografiat, permit utilizatorului să aleagă între două moduri majore de stocare a imaginii: fie imaginea este stocată în spațiul sRGB lăsând dispozitivul să aleagă metoda de compresie a gamei dinamice, fie este stocată sub formă brută, urmând ca toți pașii asociați cu compresia gamei dinamice să fie efectuați ulterior de către fotograf. Această a doua variantă de stocare produce imagini în format brut, așa-numitul ***RAW format***.

Pe lângă informația brută provenită de la senzori, formatul RAW conține informații conexe despre imagine (informații care sunt stocate de altfel și în imagini JPEG), plus informații legate de senzor și o reprezentare redusă a imaginii pentru identificare ei vizuală, fără a fi necesară procesare suplimentară. RAW este un format relevant în fotografia digitală, alături de Tiff sau JPEG.

RAW este o alternativă superioară față de TIFF. Fișierele RAW sunt mai reduse și pot reține chiar mai multe informații despre imagine .

Fisierele RAW pot oferi o plajă mult mai extinsă de ajustări decât fișierele JPEG: peste 50% mai multe detalii la expunere (acest lucru depinde însă de capacitatea senzorului). Doarece fișierele RAW conțin o gamă mai variată de culori, balansul pe alb poate fi corectat fără pierderi mari de detalii. Folosind convertoarele RAW cum ar fi plugin-ul Photoshop Camera, RAW este mai ușor și produce rezultate mult mai bune.

Marele avantaj al lucrului în formatul RAW este dat de faptul că asigură o mare flexibilitate, fiind formatul care oferă și cel mai mult spațiu pentru modificarea imaginii, păstrând în același timp calitatea.

Formatul RAW oferă :

- Imagini cu cel mai mare interval dinamic - putând înregistra mai bine zonele luminoase și umbrele
- Imagini al căror balans de alb poate fi reglat în orice moment
- Un mai bun control al clarității în comparație cu alte formate de fișiere
- O mai mare flexibilitate în corectarea erorilor de expunere

RAW-urile conțin toate informațiile înregistrate de cameră în momentul captării imaginii, însă au două mari defecte: ocupă mult spațiu și necesită un editor specializat pentru a le deschide.

2.1.3. Formate audio

AIFF Audio Interchange File Format

Audio Interchange File Format Audio File Interchange Format (AIFF) este un format audio standard utilizat pentru stocarea datelor de sunet pentru calculatoare personale și alte dispozitive electronice audio. Formatul a fost co-dezvoltat de Apple Computer, în 1988 pe baza Interchange Electronic Arts - Format de fișier (IFF, utilizate pe scară largă la sistemele Amiga) și este cel mai frecvent utilizat pe sisteme de computere Apple Macintosh.

Datele audio într-un fișier AIFF standard sunt impulsuri-modulare necomprimate cod (PCM). Există, de asemenea, o variantă comprimată de AIFF cunoscut sub numele de AIFF-C sau aifc, cu diverse codec-uri de compresie definite.

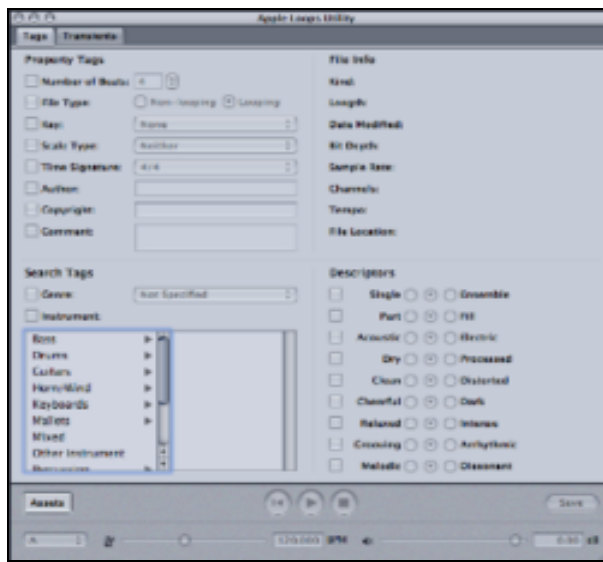


Fig. 2.4. Apple Loops Utility

Standard AIFF este un format de conducere (împreună cu SDII și WAV), utilizat la nivel profesional în aplicații audio și video, iar spre deosebire de formatul mai bine cunoscut lossy MP3, acesta este non-comprimat (care SIDA rapidă de streaming de mai multe fișiere audio de la disc pentru a cererii), și fără pierderi. Ca orice format non-comprimat fără pierderi, folosește spațiu pe disc mult mai mult decât MP3-aproximativ 10 MB pentru un minut de audio stereo la o rată eșantion de 44,1 kHz și un eșantion de 16 biți. În plus față de datele audio, AIFF poate include date de tip bucla punct, iar nota muzicala a unei probe, pentru a fi utilizate de către prelevatorii hardware și aplicații muzicale.

Extensiile pentru formatul standard AIFF sunt .aiff sau .aif. Pentru variantele comprimate se presupun a fi .aifc, dar și .aiff sau .aif sunt acceptate de aplicațiile audio.

Apple a mai creat recent o extensie la formatul AIFF sub forma Apple Loops (buclă) folosit mai ales de GarageBand și Logic Audio, care permit includerea de date pentru trecerea ritmului și a tempo-ului în aplicații mai comune.

Un fișier AIFF este împărțit într-un număr de părți fiecare identificată de un ID .

Tipuri de părți în fișierele AIFF:

- Common Chunk (required)
- Sound Data Chunk (required)
- Marker Chunk
- Instrument Chunk
- Comment Chunk
- Name Chunk
- Author Chunk
- Copyright Chunk
- Annotation Chunk
- Audio Recording Chunk
- MIDI Data Chunk
- Application Chunk

WAV - Waveform audio format

WAV este o prescurtare a „Waveform Audio Format”, adică „format de undă audio”, standard dezvoltat de IBM și Microsoft pentru pastrarea fișierelor audio pe PC. Este un format fără compresie folosit de obicei pentru secvențe scurte de sunet și înregistrări vocale. Raspindirea sa a scăzut în ultimul timp, datorită mai ales dimensiunii mari a fișierelor și a apariției formatului mp3.

Acest format de fișier audio stochează sunetele sub formă de unde. În funcție de mai mulți factori, un minut de sunet poate stoca o dimensiune foarte mică, precum 644 de kiloocteți sau una foarte mare, cum ar fi 27 de megaocteți.

WAV este compatibil cu Windows, Macintosh, și sisteme de operare Linux.

Deși un fișier WAV poate conține fișiere audio comprimate, forma cea mai comună WAV conține fișiere necomprimate audio în module liniare cod format (LPCM). Formatul standard de fișiere audio pentru CD-uri, de exemplu, este LPCM-codificate, care conține două canale de 44100 de eșantioane pe secundă, 16 biți per eșantion. Deoarece LPCM utilizează o metodă de stocare necomprimată, care pastrează toate mostrele de pista audio, utilizatorii profesioniști sau experți audio poate folosi formatul WAV pentru calitate maximă audio. WAV audio poate fi, de asemenea, editat și manipulat cu o relativă ușurință utilizând software-ul. Format WAV sprijină comprimatul audio, folosind, pe Windows, Audio Compression Manager.

Datorită structurii simple, bazate pe formatul IFF, este folosit cu o varietate de aplicații software.

Datorită dimensiunilor mici ale fișierelor comprimate, dar formate lossy cum ar fi MP3, ATRAC, AAC, Ogg Vorbis, WMA acestea sunt folosite pentru a stocarea și transferul datelor audio. Dimensiunile mici ale acestor fișiere permit o transmitere rapidă pe Internet precum și un consum mai mic de spațiu pe mass-media de memorie. Cu toate acestea, formatele lossy comprimate pierd din calitatea audio.

Formatul WAV este limitat la fișiere de cel mult 4GB, ceea ce ar însemna circa 6,8 ore de calitate CD-audio, motiv pentru care a fost necesară depășirea acestor limite prin crearea formatului W64.

De când frecvența transferului de date a unui fișier WAV a variat de la 1Hz la 4,3 GHz și numărul de canale este mai mare de 65536, fișierele .wav au fost folosite pentru date nonaudio.

WMA- Windows Media Audio

Windows Media Audio este o tehnologie audio de compresie de date dezvoltată de Microsoft. Numele poate fi utilizat pentru a se referi la formatul de fișier audio sau codecurile sale audio. Ea face parte din cadrul Windows media.

WMA este un format „proprietate” de compresie audio ce oferă posibilitatea protejării fișierelor audio împotriva copierii ilegale printr-o tehnică numită „gestiune numerică de drepturi”.

WMA există în patru forme (codecuri):

- WMA standard- codec original un concurent pentru MP3 și codec-uri Real Audio, cel mai răspândit pe Internet fiind „lizibil” pentru mai multe browser-e;

- WMA Pro- codec mai nou și mai avansat, susține mai multe canale și are o înaltă rezoluție audio, dar mai puțin răspândit;

- WMA Lossless- comprimă datele audio, fără pierderi de fidelitate audio, oferind o calitate sonoră identică cu a originalului,

- WMA Voice – în special folosit pentru codificarea vocii, comprimarea se face cu o viteză mai mică sau egală cu 20kbit/s.

Primul codec WMA s-a bazat pe munca lui Henrique Malvar și a echipei sale care au lucrat la proiectul MSAudio. Primul codec finalizat a fost prevăzut ca o continuare MSAudio 4.0. Mai târziu a fost lansat oficial ca Windows Media Audio, ca partener a Windows Media Technologies 4.0.

Windows Media Audio Standard (WMA) este cel mai utilizat codec din cele enumerate mai sus.

Versiuni disponibile WMA: Windows media Audio 2 în 1999, Windows Media Audio 7 în 2000, Windows Media Audio 8 în 2001 și Windows Media Audio 9 în 2003 (cu codecurile: Windows Media Audio 9 Professional, Windows Media Audio 9 Lossless și Windows Media Audio 9 Voice) , Windows Media Audio Professional 10 în 2007.

Windows Media Audio 9 are o frecvență de 44,10 sau 48 kHz, ocupă 16 B, similar CD-ului, sunetul având aceeași calitate și o rată cuprinsă între 64 și 192 kB/s. În această variantă calitatea sunetului este cu 20% mai bună decât a versiunii din 2003. Comparativ cu MP3 calitatea sunetului este mai mare la aceeași rată de biți.

AVI Audio Video Interleave

AVI este un format multimedia de tip „container” conceput ca parte a Video pentru Windows de către Microsoft în 1992 pentru a stoca date audio-video, permițând sincronizarea audio-video.

Cele mai multe formate AVI folosesc extensiile dezvoltate de Matrox Open DML Group, care neoficial se numesc AVI 2.0.

Într-un fișier AVI, fiecare element audio sau video poate fi comprimat de orice codec. Formatul DivX este codec-ul folosit pentru datele video și formatul Mp3 pentru datele audio dar

pot fi folosite și alte codec-uri, cum ar fi: XviD, MPEG pentru video și Mp2 și WAV pentru audio.

În acest format pot fi reunite într-un singur fișier a unei piste video și până la 99 canale audio, ceea ce permite de exemplu subtitrarea unui film în mai multe limbi.

Multe fișiere AVI folosesc acum extensii ale formatelor dezvoltate de Matrox Open DML în februarie 1996. Aceste fișiere sunt acceptate de Microsoft și sunt numite neoficial AVI 2.0.

DAT- Digital Audio Tape

Digital Audio Tape (DAT sau R-DAT) este un suport de înregistrare și mijloc de redare audio numeric pe bandă magnetică de 3,14 mm dezvoltat de Sony în anii 80. Ca aspect este similară casetei audio dar este mai mică 73x54x10,5mm și face parte din categoria mijloacelor de înregistrare digitale asigurând copierea exactă a unui alt material digital, spre deosebire de copierea pe CD.

Tehnologia DAT este stâns legată de cea a magnetoscoapelor, folosind un cap rotativ pentru înregistrare diferența constând în cantitatea mult mai mare de date ce poate fi stocată.

DAT nu a fost prima casetă audio- digitală, a mai existat o tentativă în 1972 în Japonia, dezvoltată de Denon pentru înlocuirea discurilor de vinilin, dar nu s-a dezvoltat într-un produs de consum.

În 1976 primul succes comercial al DAT-ului a fost dezvoltat de Sony folosind un aparat de înregistrare fabricat de Honeywell conectat la un hardware ce folosea programe de codare- decodare create de Sony.

La scurt timp după Sony, 3M a introdus propria linie de înregistratoare și formate ale acestor casete pentru a putea fi folosite în studiourile de înregistrări din Minneapolis, Minnesota.

La sfârșitul anilor 80, Recording Industry Association of America a făcut lobby împotriva introducerii dispozitivelor DAT în SUA, pentru a preveni folosirea acestora pentru copierea discurilor sau casetelor.

Acest gen de acțiuni au dus la încetarea fabricării Replayer-ului DAT pe scară largă.

Acest gen de format de stocare a fost utilizat pe scară largă în industria de înregistrări audio profesionale și în radiodifuziune datorită faptului că permitea codificarea fără pierderi și cantitatea de date stocate varia între 1,3 și 80 GB pe o bandă de 60 – 180 metri.

Transferul de date de pe bandă pe hard disc se face prin intermediul unei conexiuni digitale de tip AES/EBU.

WMV- Windows Media Video

WMA este numele unei familii de codec-uri video dezvoltate de Microsoft. Pe Internet se întâlnește frecvent acest tip de fișier. Ca cea mai mare parte a codec-urilor „proprietate” această tehnologie nu corespunde cu cea a normelor internaționale (ISO/IEC MPEG), Microsoft dorind intrarea acestor codec-uri pe cele mai multe obiecte de larg consum (telefoane, DVD playere, decodoare TV) altele decât cele obișnuite, a încercat introducerea versiunii 9, WMV9, în consorțiile industriale (3GPP, DVB, ATSC, DVD-Forum, Blu-ray).

Fiind în parte în neconcordanță cu normele MPEG-4 AVC/H 264, ce au fost inițiate și dezvoltate de crecetători din companii internaționale și universități din toată lumea, pe de altă parte fiind obligați de consorțiile industriale de a prezenta un codec standardizat, Microsoft a înscris codec-ul WMV9 la SMPTE (society of Motion Picture and Television Engineers) pentru a-l standardiza sub numele VC-1.

În 2005, codec-ul VC-1 devine stabil și are statut de standard internațional.

VC-1 are versiune high definition, a fost implementat pe HD DVD de la NEC și pe discurile Blu-ray de la sony, care probabbil vor fi înlocuitorii DVD- playerelor de astăzi.

Codecuri Windows Media Video :

- Windows Media Video v.7
- Windows Media Video Screen v.7
- Windows Media Video v.8
- Windows Media Video v.9
- Windows Media Video v.9 Screen
- Windows Media Video 9 Advanced Profile (FourCC: WVC1)
- Windows Media Video 9 Image
- Windows Media Video 9.1 Image
- VC-1

2.1.4. Formate video

Standardele MPEG. MPEG este acronimul pentru Moving Picture Experts Group (grupul experților în filme) creat de organizația internațională pentru standardizare (ISO).

Grupul a fost format pentru a stabili standarde de compresie și transmisie audio și video. Prima întâlnire a grupului a avut loc în orașul Ottawa din Canada în mai 1988.

Metodologia de compresie MPEG este considerată asimetrică deoarece codorul este mai complex decât decodorul. Codorul trebuie să fie algoritmic sau adaptabil pe când decodorul execută acțiuni fixe. Acest fapt este considerat un avantaj în aplicațiile de difuzare în masă unde numărul de codoare scumpe este mic dar numărul decodoarelor ieftine este mare.

Standardele MPEG sunt compuse din mai multe părți. Fiecare parte acoperă un anumit aspect al întregii specificații. Standardele specifică de asemenea și profilele și nivelele. Profilele sunt destinate să definească un set de instrumente care sunt valabile, iar nivelele definesc gama de valori potrivite pentru proprietățile acestora. MPEG a standardizat următoarele formate de compresie și standarde auxiliare:

- MPEG-1: Acest standard a apărut în anul 1993. Utilizat la codarea filmului și a sunetului asociat pentru medii digitale de stocare cu viteze până la 1,5 Mbit/s. Este primul standard pentru compresie audio și video. Principiul de proiectare a fost permiterea codării filmului și a sunetului asociat la rata de biți a unui disc compact (CD). Pentru aceasta MPEG-1 reduce rata semnalului de transmisie pentru imagini, dar utilizează și o frecvență a pozelor de doar 24-30 Hz, rezultând o calitate moderată. Include popularul format de compresie audio MP3.

- MPEG-2: Standardul a fost lansat în anul 1995. Era utilizat pentru codarea generică a filmelor și a informațiilor audio asociate. MPEG-2 suportă întrepătrundere și o înaltă definiție. Standardul este considerat important deoarece a fost ales ca și schemă de compresie pentru televiziunile digitale ATSC, DVB și ISDB, servicii digitale prin satelit precum Dish Network, SVCD, DVD și Blu-ray.

- MPEG-3: Se ocupă de standardizarea compresării scalabile și de rezoluții multiple și a fost creat pentru compresie HDTV însă a fost descoperit că este redundant și a fost contopit cu MPEG-2, astfel că nu există standardul MPEG-3. A nu se confunda MPEG-3 cu MP3 care este stratul audio 3 al standardului MPEG-1.

- MPEG-4: A apărut în 1998. A fost creat pentru codarea obiectelor audio-vizuale. MPEG-4 folosește unelte de codare cu complexitate sporită pentru a atinge factori de compresie mai mari decât MPEG-2. Pe lângă codarea video mai eficientă, MPEG-4 se apropie mai mult de aplicațiile grafice. În profile mai complexe, decodorul MPEG-4

devine un eficient procesor de interpretare grafică, iar fluxul de date compresate descrie forme tri-dimensionale și texturi ale suprafețelor.

Formatul MOV

Formatul Quicktime (.mov) funcționează ca un recipient care conține una sau mai multe piste, fiecare conținând un singur tip de date: audio, video, efecte sau text. Fiecare pistă conține un flux de date codate digital (folosind un codec specific) sau o referință către fluxul de date aflat într-un alt fișier. Pistele sunt menținute într-o structură de date ierarhizată care este formată din obiecte numite atomi. Un atom poate fi părintele altor atomi sau poate conține date media, dar nu poate fi ambele.

Abilitatea de a conține referințe abstracte de date pentru datele media și separarea datelor media de listele de editare a pistelor duc la faptul că formatul Quicktime este printre cele mai potrivite pentru editare, fiind capabil să importe și să modifice datele pe loc fără necesitatea rescrierii întregului fișier.

Formatul OGG

OGG este un format standard gratuit, deschis dezvoltării publice, întreținut de fundația Xiph.Org.

Creatorii formatului susțin că acesta nu este restricționat de brevete de software și este proiectat să producă un eficient flux de date și o manipulare eficientă a datelor multimedia digitale de calitate înaltă.

Numele „Ogg” se referă la formatul care poate multiplexa mai multe codecuri gratuite, deschise dezvoltării publice, pentru date audio, video, text și meta.

În cadrul arhitecturii multimedia Ogg, Theora furnizează stratul video, în timp ce codecul orientat spre muzică Vorbis acționează ca strat audio. Ca și strat audio mai pot fi utilizate și Speex, FLAC sau OggPCM.

Termenul „Ogg” este deseori folosit ca referință la formatul audio Ogg Vorbis care este de fapt audio codat folosind Vorbis în formatul de stocare Ogg. Anterior, extensia .ogg era folosită pentru orice conținut distribuit prin Ogg, însă începând cu anul 2007, fundația Xiph.Org a cerut ca extensia .ogg să fie utilizată doar pentru codecul Vorbis pe baza îngrijorărilor legate de compatibilitate inversă. Fundația a decis să creeze un nou set de extensii ale fișierelor și tipurilor media pentru a descrie diferitele tipuri de

conținut precum .oga pentru fișiere cu conținut exclusiv audio, .ogv pentru video cu sau fără audio (inclusiv Theora), și .ogx pentru aplicații.

Formatul AVI

AVI, acronimul de la „Audio Video Interleave” (Interclasare audio video), este un format de stocare multimedia digital introdus de Microsoft în 1992 ca parte a tehnologiei „Video pentru Windows”. Fișierele AVI pot conține atât date audio cât și date video într-un fișier care permite redare sincronizată audio-video, însă nu e necesar ca acestea să apară împreună (exemplul unui clip fără sunet).

AVI este o derivare a formatului intitulat RIFF care împarte datele unui fișier în bucăți. Fiecare bucată este identificată de o etichetă FourCC(cod de patru caractere) . Un fișier AVI ia forma unui fișier de tip RIFF format dintr-o singură bucată care este divizată în două părți mai mici obligatorii și o parte opțională.

Prima parte obligatorie este identificată de eticheta „hdrl” și reprezintă antetul fișierului care conține informații legate de video precum înălțimea, lățimea și rata cadrelor. A doua parte obligatorie este identificată de eticheta „movi” și conține toate datele audio-video care alcătuiesc filmul. Partea opțională este identificată de eticheta „idx1” și indexează adresele relative ale bucăților de date din fișier.

Asemănător formatului RIFF, datele audio-video conținute în partea „movi” pot fi codate sau decodate de un software numit codec. La crearea fișierului, codecul translatează între datele neprelucrate și tipul de date compresate folosit în fișier. Un fișier AVI poate stoca datele audiovideo în aproape toate formele de compresie, inclusiv : Full Frame (necompresat), Intel Real Time (Indeo), Cinepak, Motion JPEG, Editable MPEG, VDOWave, ClearVideo/RealVideo, QPEG și MPEG-4 video.

Secțiunea ce urmează nu descrie formatul OpenDML AVI (așa numit-ul "AVI 2.0"), dezvoltat de compania Matrox, ci doar specificațiile Microsoft.

Fișierele AVI sunt identificate după FourCC-ul „AVI ” (inclusiv spațiul, al patrulea caracter) din antetul RIFF. Toate AVI-urile au cel puțin două liste esențiale de părțicele, care conțin formatele fluxurilor, respectiv fluxurile de date propriu-zise. Un AVI poate avea, opțional, o porție ce conține locațiile celorlalte porții de date din fișier. Atunci AVI-ul ar avea următoarea structură:

RIFF ('AVI '

```
LIST ('hdrl' ... )
LIST ('movi' ... )
['idx1' (<AVI Index> ) ]
```

Aici, „hdrl” și „movi” sunt ID-urile celor două liste esențiale despre care vorbeam. „hdrl” definește formatul de date și este prima listă necesară. „movi” conține datele pentru secvența audiovideo și este cea de-a doua listă necesară. Porția cu ID-ul „idx1” conține indecșii. Fișierele TREBUIE să aibă aceste trei componente în ordine. Listele „hdrl” și „movi” au ca date subporții (subchunks). Următorul exemplu prezintă într-un mod ceva mai extins structura unui fișier RIFF AVI:

```
RIFF ('AVI '
LIST ('hdrl'
'avih'(<Main AVI Header>)
LIST ('strl'
'strh'(<Stream header>)
'strf'(<Stream format>)
[ 'strd'(<Additional header data> ) ]
[ 'strn'(<Stream name> ) ]
LIST ('movi'
{SubChunk | LIST ('rec '
SubChunk1
SubChunk2
['idx1' (<AVI Index> ) ]
```

Antetul AVI principal

Lista „hdrl” începe cu sub-chunkul „avih” (AVI header), ce conține antetul AVI principal, informații valabile pentru întregul fișier AVI, cum ar fi numărul de fluxuri din fișier (en. streams –spre exemplu, un fișier și video și audio are două fluxuri), lățimea și înălțimea secvenței AVI (în pixeli), numărul de microsecunde dintre cadre, numărul total de cadre și alte informații. Pentru mai multe detalii, rămărește linkul [http://msdn.microsoft.com/en-us/library/ms779632\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/ms779632(VS.85).aspx), unde este definită structura AVIMAINHEADER.

Antetele fluxurilor AVI

După „avih” (AVI main header) urmează una sau mai multe liste „strl”. Pentru fiecare flux de date (audio, video) este necesară o listă „strl”. O astfel de listă conține informații despre fluxul asociat ei și trebuie să aibă o porție pentru antetul fluxului („strh” – stream header chunk) și o porțiune pentru formatul fluxului („strf” – stream format chunk). Opțional, o listă „strl” poate conține porțiunile „strd” (stream header data chunk) și „strn” (stream name chunk).

Porțiunea cu ID-ul „strh” alcătuiește o structură AVISTREAMHEADER. Pentru mai multe informații, urmează linkul [http://msdn.microsoft.com/en-us/library/ms779638\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/ms779638(VS.85).aspx).

Porțiunea cu ID-ul „strf” trebuie să urmeze imediat după „strh”. Acesta descrie formatul datelor din flux, deci depinde de tipul acestuia. Pentru fluxuri video, informația regăsită în porția „strf” alcătuiește o structură BITMAPINFO (vezi [http://msdn.microsoft.com/en-us/library/dd183375\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/dd183375(VS.85).aspx)), iar pentru fluxuri audio, o structură WAVEFORMATEX ([http://msdn.microsoft.com/en-us/library/ms788112\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/ms788112(VS.85).aspx)).

Dacă porțiunea „strd” este prezentă, atunci ea trebuie să urmeze porțiunii „strf”. Formatul și conținutul ei sunt definite de driver-ul codec. În general, driverele folosesc informația din această porțiune pentru configurări. Aplicațiile care redau, creează sau modifică fișiere AVI nu au nevoie să interpreteze informația. Ele doar o transferă la sau de la driver ca un simplu bloc de memorie.

Porțiunea opțională cu ID-ul „strn” conține un șir de caractere (cu caracterul nul la sfârșit) descriind fluxul.

Anteturile fluxurilor din lista „strl” sunt asociate cu fluxurile propriu-zise (datele audio sau video) din lista „movi” în funcție de ordinea în care apar listele „strl”. Astfel, prima listă „strl” se aplică fluxului 0, a doua fluxului 1, și așa mai departe.

Datele fluxurilor (lista „movi”)

După toate informațiile din antet, urmează lista „movi”, ce conține datele propriu-zise, altfel spus, cadrele pentru fluxurile video și sample-urile pentru fluxurile audio. Porțiunile de date pot fi așezate direct în conținutul listei „movi” sau grupate în liste „rec” ce sunt așezate, bineînțeles în interiorul lui „movi”. Listele-grup cu ID-ul „rec” sunt

folosite pentru a determina aplicația de redare (orice AVI Player) să citească datele din fișier dintr-o singură parcurgere. Acest lucru este util, spre exemplu, pentru vizionarea filmelor în format AVI de pe un CD.

În „rec” se află, așadar, o mulțime de subporțiuni, care conțin date audio, sau vizuale. ID-ul fiecărei sub-porții este un FourCC format din doi octeți ce eprezintă indexul fluxului de date de care aparține (se ține cont de ordinea în care apar listele „strl” în fișier) la care se adaugă două caractere ASCII ce definesc tipul de date din porțiune(deci în total $2+2=4$ octeți).

Porția AVI Index („idx1”)

Această porțiune este opțională, putând urma după lista „movi”. Conține o listă cu ID-urile porțiunilor de date (cadrele, sample-urile audio), mărimea, locațiile lor în fișier. În listă sunt incluse, de asemenea, și grupurile de porțiuni de date, porțiunile „rec”. Porțiunea „idx1” alcătuiește o structură AVIOLDINDEX ([http://msdn.microsoft.com/en-us/library/ms779634\(VS.85\).aspx](http://msdn.microsoft.com/en-us/library/ms779634(VS.85).aspx)). Dacă fișierul AVI conține această porțiune finală, atunci câmpul dwFlags din structura AVIMAINHEADER (adică porția „avih”) trebuie să aibă setată caracteristica AVIF_HASINDEX.

Alte părți pot fi cele cu ID-ul „JUNK”. Acestea sunt folosite pentru alinierea datelor în fișierul AVI.

O aplicație ar trebui să ignore conținutul acestei părți.

Digital Video Data (DV DATA) în fișierele AVI

În cele ce urmează va fi prezentată specificația Microsoft pentru stocarea digitală a secvențelor video. Conformându-se la această specificație, fișierele AVI sunt și vor fi compatibile cu versiuni curente și viitoare ale arhitecturii video pentru platforma Windows.

În cele ce urmează, se descrie structura AVI-urilor ce conțin DV data. Se definesc FourCC-uri specifice pentru fluxurile DV data și obiecte de comandă (en. handlers) pentru compresia / decompresia fluxurilor. De asemenea, se definește structura unui flux de DV data și se specifică cele două metode de stocare a DV data într-un fișier AVI.

2.1.5. Formate multimedia

Multimedia (MM) înseamnă abilitatea de a achiziționa, manipula, combina și reda informații de la o mare varietate de medii, ce includ text, grafică, sunet, imagine fixă sau video. Multimedia nu este o tehnologie, ci mai degrabă un termen ce descrie un număr de tehnologii care lucrează împreună. Multimedia a transformat interacțiunea om-calculator și a permis crearea unei noi familii de produse în diverse domenii precum:

- *acces la cunoaștere* – multimedia este modul cel mai rapid, eficient și ieftin de a permite indivizilor accesul la informare, punându-le la dispoziție adevărate enciclopedii electronice;
- *administrarea documentelor și a înregistrărilor* – datorată apariției unor documente tot mai complexe ca și conținut care trebuiau gestionate de întreprinderi și instituții comerciale;
- *educație și instruire* – furnizând materialul potrivit pentru școlarizarea și instruirea unor categorii variate de subiecți;
- *reclame* – în mod practic nu există nici o limită în folosirea informației multimedia în astfel de aplicații;
- *controlul și monitorizarea proceselor în timp real* – împreună cu bazele de date, prezentările multimedia de informații au un rol efectiv în operațiile de monitorizare și control a sistemelor precum sistemele de transport, sistemele de supraveghere a pacienților, etc.

Această nouă generație de produse permite atât integrarea elementelor multimedia în mediile existente cât și o nouă abordare în ceea ce privește procesul muncii. Așadar, multimedia poate extinde aplicațiile existente, dar poate conduce și la regândirea în mod revoluționar a procesării informației în diverse domenii cum ar fi economia, știința, arta, educația, ingineria. Folosirea multimedia generează beneficii pentru toți utilizatorii sistemelor informatice. Ea îmbunătățește calitatea și cantitatea informației prezentate utilizatorului, precum și interacțiunea om-calculator.

Organizarea documentelor multimedia

Organizarea eficientă a informației de tip multimedia cade în sarcina sistemelor de gestiune a bazelor de date multimedia (SGBDMM). Un astfel de sistem este un sistem de gestiune a bazelor de date performant care suportă tipuri de date multimedia pe lângă cele clasice, alfanumerice, și care are capacitatea de a manipula cantități foarte mari de informație multimedia integrând strâns trei tehnologii fundamentale, și anume:

- bazele de date;
- sistemele de regăsire a informației;
- sistemele de stocare ierarhică a informației.

Bazele de date multimedia pot fi definite ca sisteme de baze de date care pot stoca, manipula și interoga informație prezentată sub formă de text, audio, video, obiecte grafice și imagini statice alb-negru sau color.

Bazele de date multimedia sunt tot mai de actualitate în lumea de azi a calculatoarelor datorită faptului că oferă posibilitatea de a prelucra într-o manieră convenabilă și flexibilă diferite tipuri de obiecte cu care interacționăm în viața noastră de zi cu zi. Ca atare, într-o bază de date multimedia putem întâlni pe lângă tipurile de date clasice dintr-o bază de date și tipuri noi precum:

- *Date imagini* – acestea sunt foarte comune în bazele de date multimedia și ele pot fi simple figuri, icoane, imagini medicale precum razele X, etc.;
- *Date video* – acestea sunt date similare fișierelor video (film) care sunt stocate în baza de date;
- *Date audio* – acestea sunt date asemănătoare celor stocate în fișierele audio și de cele mai multe ori sunt de fapt melodii, voce, etc.;
- *Date document* – acestea sunt date de forma fișierelor text tradiționale în care informația este stocată sub formă de text, numai că mărimea acestor texte este foarte mare.

Obiectele multimedia sunt diferite de textul tradițional sau de documentele numerice în sensul că acestea necesită de regulă o mai mare cantitate de memorie internă și externă pentru memorare. De asemenea, operațiile aplicate obiectelor multimedia sunt

diferite (de ex. afișarea unei figuri sau vizionarea unui video clip sunt diferite de afișarea unui paragraf de text).

Un sistem de gestiune a bazelor de date multimedia (SGBDMM) trebuie să ofere un mediu adecvat pentru utilizarea și gestiunea datelor multimedia. În completarea funcțiilor tradiționale ale sistemelor de gestiune a bazelor de date, un sistem de gestiune a bazelor de date multimedia trebuie să includă funcții specifice precum:

- gestiunea unor date de tip multimedia precum imagini, video, grafică, audio, etc.;
- gestiunea unor volume mari de date multimedia;
- furnizarea unor scheme de gestiune a spațiilor de stocare performante și rentabile ca și cost;
- stocarea eficientă și gestiunea livrării datelor multimedia;
- soluții eficiente pentru indexarea și regăsirea datelor, folosind date multimedia drept criteriu de căutare;
- să suporte diferite formate pentru datele multimedia;
- să suporte diferite funcții specifice bazelor de date precum inserările, ștergerile, căutările și modificările;
- să optimizeze realizarea interogărilor și prelucrărilor.

Obiectele multimedia sunt, de regulă, obiecte binare mari BLOB (binary large objects). Este obișnuit ca un video clip să ocupe mai mult de 100 MB spațiu de stocare. Pe un sever video, este posibil ca sute de video clipuri să fie stocate. Datorită cantității uriașe de spațiu de stocare necesar, un sistem de gestiune a bazelor de date multimedia necesită un mecanism sofisticat de gestiune a spațiului de stocare, care de asemenea trebuie să aibă un preț rentabil. Schemele de gestiune a spațiului de stocare trebuie să suporte, de asemenea, și operațiile fundamentale necesare la nivelul bazei de date.

Adesea, un SGBDMM trebuie să ia în considerare și următoarele aspecte:

- compunerea și descompunerea obiectelor multimedia;
- operații ale obiectelor multimedia care implică sincronizarea acestora;

- persistența obiectelor;
- regăsirea informațiilor multimedia pe baza conținutului;
- acces concurent și mecanisme de blocare pentru procesarea distribuită;
- securitatea datelor;
- consistența datelor, integritatea referențială și refacerea datelor în caz de accident;
- tranzacții lungi și tranzacții imbricate;
- indexarea și clusterarea datelor.

Spre deosebire de SGBD-urile tradiționale, în SGBDMM replicarea datelor nu este încurajată din cauza volumului mare de date vehiculat.

Pentru aplicațiile multimedia simple, utilizarea modelului client-server pentru accesarea bazei de date multimedia, este considerată adecvată.

Aplicațiile multimedia complexe pot impune existența unui server video și utilizarea unui SGBDMM cu o arhitectură dinamică.

Arhitectura unui SGBDMM

O bază de date multimedia are arhitectura structurată pe trei nivele:

- nivelul extern al interfeței cu utilizatorii;
- nivelul conceptual de compunere a obiectelor;
- nivelul intern al mediului de stocare a datelor.

Sarcinile specifice nivelului extern al interfeței includ în principal afișarea obiectelor, procesarea interogărilor și compunerea și descompunerea obiectelor. Afișarea obiectelor permite utilizatorilor să găsească entități multimedia existente care apoi pot fi refolosite. Prin intermediul interogărilor, fie bazate pe text fie pe vizualizare, utilizatorii specifică un număr de condiții pentru entitățile multimedia dorite și obțin o listă de obiecte posibile a satisface aceste condiții. Obiectele potrivite sunt apoi reutilizate. Resursele multimedia, spre deosebire de cele de tip text sau numerice, nu pot fi localizate efectiv utilizând un limbaj de interogare bazat doar pe text. Chiar și folosind limbajul

natural pentru formularea interogării este greu de regăsit cu precizie o fotografie sau o secvență video cu un anumit conținut. Cercetările privind regăsirea informațiilor pe baza conținutului se concentrează asupra mecanismelor care permit utilizatorilor să găsească efectiv obiecte multimedia reutilizabile, incluzând poze, sunete, video și alte forme. După găsirea cu succes a obiectelor dorite, interfața bazei de date trebuie să-l ajute pe utilizator să compună/descompună documentele multimedia.

Cel de-al doilea nivel lucrează în legătură cu nivelul interfeței pentru a gestiona obiectele. În mod tipic, compunerea obiectelor necesită o serie de legături, precum legăturile de asociere, legăturile de similaritate și legăturile de moștenire pentru a specifica diferitele relații dintre obiecte. Aceste legături sunt specificate fie prin intermediul interfeței grafice utilizator a bazei de date, fie prin intermediul unor funcții ale API.

Ultimul nivel, nivelul intern al gestiunii mediului de stocare a datelor, include două aspecte legate de performanță: clusterarea și indexarea datelor. Clusterarea se referă la organizarea fizică a informației multimedia pe hard disc sau pe alt mediu optic de stocare, astfel încât atunci când se încearcă regăsirea acestor date, sistemul să fie capabil să acceseze datele binare în mod eficient. De obicei, performanța regăsirii datelor trebuie să garanteze un anumit nivel de calitate a serviciului oferit (QoS) și să realizeze sincronizarea obiectelor multimedia. Indexarea înseamnă că un mecanism de localizare mai rapid este esențial pentru a găsi adresa fizică a unui obiect multimedia. Câteodată, schemele implică date complexe sau structuri de fișiere.

Modele de date

Datele multimedia de tip continuu, precum video și audio, implică folosirea conceptelor specifice precum fluxul datelor, compoziție temporală, eșalonare în timp și sincronizare. Aceste concepte sunt diferite de cele utilizate în cazul modelelor convenționale de date și în consecință modelele de date convenționale nu sunt compatibile cu sistemele de baze de date multimedia. O problemă cheie cu care se confruntă un SGBDMM este descrierea structurii mediilor dependente de timp într-o formă adecvată interogării, actualizării, regăsirii și prezentării datelor.

În acest sens, printre primele realizări se găsește *modelul VideoSTAR* care reprezintă un standard pentru baze de date video și este proiectat ca suport pentru aplicații video care partajează și reutilizează date video și metadate. Metadatele sunt implementate ca o bază de date complexă. Sistemul conține un model principal de indexare valabil pentru multe tipuri de aplicații cu câteva entități predefinite precum: persoane, evenimente etc. Fiecare aplicație poate îmbunătăți puterea descriptivă a modelului de bază cu tipuri specifice de indecși, idee folosită în multe alte modele.

Un alt model de date descris de Zhong prezintă o schemă generală pentru modelarea obiectelor video, care încorporează caracteristici vizuale de nivel scăzut și grupări ierarhice. Modelul oferă o structură generală pentru extragerea obiectelor video, pentru indexare, împărțirea pe categorii și prezentarea noilor segmente video și algoritmi de urmărire a obiectelor pe baza culorii relevante și caracteristicilor traiectoriei. Prin obiecte video, autorii se referă la obiectele de interes, incluzând cele mai importante regiuni ale imaginilor de nivel scăzut, obiectele în mișcare și grupe de obiecte care satisfac constrângeri spațio-temporale. Aceste obiecte video pot fi extrase automat prin segmentarea obiectelor și utilizarea unor mecanisme de urmărire a obiectelor. După aceea obiectele sunt stocate într-o bibliotecă digitală. Obiectele pot fi grupate și fiecărui grup îi pot fi asociate adnotări semantice de nivel ridicat. Se propune un model ierarhic de reprezentare a obiectelor în care obiecte din niveluri diferite pot fi indexate, căutate și grupate în concepte de nivel ridicat. Interogările bazate pe conținut pot doar regăsi segmente video și pot seta punctele de intrare în video în funcție de caracteristicile spațio-temporale ale obiectelor video interogate și nu se referă la regăsirea unităților video dedicate din întregul flux video.

Cercetările mai recente se concentrează pe sisteme informaționale multimedia, în care SGBDMM manipulează datele multimedia într-o manieră structurată și care folosesc sisteme de recunoaștere a datelor multimedia.

Multe din modelele conceptuale ale datelor multimedia sunt implementate ca scheme de baze de date sau ca sisteme obiectuale sau relaționale și folosesc serviciile standard ale SGBD-urilor precum independența datelor (abstractizarea datelor), neutralitatea aplicației (deschisă), controlarea accesului multiutilizator (control concurrent), toleranța la erori (tranzacții, recuperare) și controlul accesului.

Aplicațiile multimedia sunt, în general, distribuite folosind mai multe servere pentru a satisface cerințele de stocare și prelucrare. Pentru gestiunea distribuției eficiente a fluxului datelor acestea pot fi gestionate de SGBD-uri sub forma BLOB-urilor (*Binary Large Object*) și pentru gestiunea accesului la distanță, folosind diferite protocoale de comunicare, diverse tehnologii de acces la date: ODBC (*Open Database Connectivity*), ADO (*ActiveX Data Object*), OLEDB (*Object Linking and Embedding DataBase*), JDBC (*Java DataBase Connectivity*).

Proiectarea de modele pentru schimbul de date, adică descrierea datelor împreună cu descrierea fluxului datelor sau ca un flux de sine stătător este un punct important într-un sistem multimedia distribuit. Datele multimedia descriptive pot oferi, de exemplu, componente de rețea active cu date utile în administrarea și îmbunătățirea politicilor de scalare a mediilor și în construirea de buffere cu conținut. Datele descriptive pot asista clientul în selectarea modului de livrare a datelor multimedia, de exemplu vizualizare imediată sau nu a datelor.

Standardul SQL/MM și MPEG-7

Organizația internațională *International Organization for Standardization and International Electrotechnical Commission* (ISO-IEC), grupul de lucru SQL a dezvoltat o extensie a limbajului SQL numită SQL/MM (Multimedia).

Pe de altă parte, standardul MPEG-7 conține o modalitate de descriere a conținutului diferitelor date multimedia. Elementele MPEG-7 sunt codificate ca documente XML pe baza regulilor limbajului pentru descrierea definițiilor datelor (DDL). DDL oferă mijloacele pentru definirea structurii, conținutului și semanticii documentelor XML. Atât MPEG-7 cât și SQL/MM introduc câte un model conceptual pentru datele multimedia, în vederea utilizării în sistemele de baze de date multimedia.

De exemplu, definirea tipului folosit pentru imagini statice în cele două limbaje se face cu ajutorul lui *SI_StillImageType* în SQL/MM și a tipului *StillRegionType* în MPEG-7.

Pe lângă definirea unei imagini sunt propuse tipuri pentru definirea caracteristicilor folosite la compararea imaginilor. De exemplu, *SI_FeatureList* oferă o listă a tuturor caracteristicilor de nivel scăzut disponibile, listă care conține în principal

mijloace de indexare a culorii și texturii. Pentru culoare este propusă o valoare-histogramă (*SI_ColorHistogram*) și două caracteristici nonhistogramă: o valoare medie (*SI_AverageColor*) și un vector al culorilor dominante (*SI_PositionalColor*). Caracteristica textură (*SI_Texture*) conține valori reprezentând caracteristicile texturii imaginii.

Culoarea medie a imaginii e definită ca media aritmetică a valorilor culorii tuturor pixelilor imaginii și este similară - dar nu e identică - cu culoarea predominantă a imaginii. Histograma culorii unei imagini este specificată ca distribuția culorilor unei imagini măsurată de-a lungul spectrului celor N culori (64 sau 128 sau 256). Textura unei imagini e măsurată folosind trei factori: asprimea, contrastul și direcționalitatea. Asprimea indică mărimea elementelor repetitive ale imaginii. Contrastul este variația luminii față de întuneric. Direcționalitatea indică existența sau nu a unei direcții predominante. Textura este folosită pentru căutarea unei imagini cu un format anume. Comparativ cu SQL/MM, modelul conceptual al MPEG-7 este mai bogat. La nivelul secvenței DDL ce definește o imagine statică sunt înglobate atât elementele descriptive ale unei imagini statice, cât și descrierile semantice și elementele de descompunere spațială ale imaginii, care nu sunt acoperite de SQL/MM.

În descrierea altor tipuri de date multimedia se păstrează aceeași abordare.

În concluzie, SQL/MM acoperă partea sintactică a descrierilor multimedia, dar nu are nici o modalitate de a descompune o imagine și nu există nici o soluție pentru utilizarea descrierilor semantice.

Majoritatea modelelor pentru datele multimedia utilizate de SGBDMM folosesc variante ale SQL/MM. De exemplu, Oracle 10g folosește în componenta *InterMedia* tipul *ORDImage* care are elemente similare cu tipul *SI_StillImage* prezentat, îmbunătățit cu informații detaliate despre formatul și compresia imaginii.

SQL/MM se integrează bine cu SGBD și de aceea oferă soluții pentru indicarea drepturilor de acces, soluții de recuperare a datelor și este operațional, în sensul că metodele pentru regăsire și prelucrarea imaginilor sunt asociate cu ierarhia tipurilor de date.

Pentru ca MPEG-7 să ofere funcționalitatea SQL/MM, acesta trebuie folosit împreună cu un SGBD care să stocheze și să permită indexarea documentelor MPEG-

7. Dezavantajul major al variantei MPEG-7 este că nu oferă mecanisme de regăsire și prelucrare proprii.

Cu toate acestea, formatul datelor folosit de SQL/MM nu permite interschimbul datelor multimedia și a metadatelor cu alte sisteme software.

Produse pentru gestiunea bazelor de date multimedia

SGBD-urile relaționale și cele obiectuale sunt, în acest moment, cele mai promițătoare platforme tehnice pentru implementarea unui model de date multimedia, pentru prelucrarea și optimizarea interogărilor, deoarece oferă următoarele facilități:

- permit definirea de noi tipuri de date plecând de la cele existente și definirea de funcții și metode de acces la datele asociate;
- oferă soluții extensibile de indexare, prelucrare și optimizare;
- gestionează datele multimedia extern și intern. Gestiunea externă a datelor multimedia presupune stocarea acestora ca fișiere de sine stătătoare, iar în baza de date sunt păstrate doar numele acestor fișiere. Gestiunea internă presupune stocarea datelor în baza de date sub forma obiectelor distincte.

SGBDMM se bazează, în principal, pe serviciile sistemului de operare pentru stocarea și regăsirea fișierelor.

Primul val al SGBDMM a fost, la mijlocul anilor '90, reprezentat de *MediaDB* (numit acum *MediaWay*), *JASMINE* și *ITASCA*.

Aceste sisteme puteau manipula diverse tipuri de date și ofereau mecanisme pentru inserarea, interogarea, regăsirea și actualizarea datelor. Majoritatea sistemelor au dispărut de pe piață după câțiva ani.

Al doilea val de SGBDMM comerciale manipulează conținutul multimedia prin intermediul tipurilor de obiecte complexe definite pentru diferite tipuri de medii. Cele mai avansate soluții sunt oferite de Oracle 10g, IBM DB2 și IBM Informix.

ORACLE interMedia

InterMedia este o componentă care extinde funcționalitățile sistemului de gestiune a bazelor de date Oracle permițând stocarea, gestiunea și regăsirea datelor

multimedia: a imaginilor, a secvențelor video, a datelor audio și a altor tipuri media eterogene, într-o manieră integrată cu tipuri de date tradiționale.

Oracle *interMedia* permite accese SQL standard utilizând servicii, operatori și metadate pentru gestiunea tipurilor de date multimedia.

Oracle *interMedia* nu controlează dispozitivele de captură multimedia și nu are funcții pentru redarea datelor multimedia ci facilitează gestiunea datele multimedia stocate în baza de date.

Conținutul multimedia poate implica rapid cantități uriașe de informație. Aplicații precum publicarea pe Web, e-comerț, gestiunea activelor media au condus la creșterea dramatică a generării și consumării de date de tip media precum poze, sunete, muzică, voce, video, etc. Pentru prima dată, securitatea, controlul administrativ, performanța, scalabilitatea și accesul gestionat profesional la sistemele de informații ale întreprinderilor este disponibil și pentru conținutul multimedia utilizat în Web site-uri corporative sau în aplicațiile multimedia.

Oracle *interMedia* utilizează obiecte de tip dată, similare cu cele din clasele JAVA sau C++ pentru a descrie date de tip imagine, audio, video. *InterMedia* permite stocarea, gestiunea și regăsirea datelor multimedia provenind din diferite surse de date. Astfel, *interMedia* gestionează datele multimedia stocate în baza de date sub forma BLOB-urilor (*binary large objects*) caz în care aceste date se află sub control tranzacțional, sau a fișierelor multimedia gestionate direct de sistemul de operare, formatul utilizat în acest caz fiind BFILE (*file-based large objects*) și a datelor multimedia stocate pe un server web. În aceste ultime două cazuri, datele multimedia sunt stocate în afara bazei de date fără a beneficia de controlul tranzacțional al acestora. În acest caz, un pointer către aceste date este stocat în baza de date sub control tranzacțional, iar datele multimedia sunt stocate în fișierul extern de tip BFILE. Datele multimedia stocate în afara bazei de date pot furniza un mecanism adecvat pentru gestiunea unor depozite media de mărimi foarte mari, pre-existente sau noi, care sunt stocate sub forma unor fișiere obișnuite pe medii de memorare read-only sau recriptibile. Aceste date pot fi importate ulterior în BLOB-uri în orice moment pentru a le putea supune controlului tranzacțional.

Metadatele, atributele și metodele obiectelor sunt întotdeauna stocate în baza de date sub controlul lui Oracle *interMedia*. Dacă datele multimedia sunt stocate în

interiorul sau în exteriorul bazei de date, *interMedia* gestionează metadatele acestora pentru toate tipurile media și automat le extrage din acestea. Aceasta face ca dezvoltatorii de aplicații să fie scutiți de povara cunoașterii caracteristicilor formatelor datelor multimedia.

InterMedia folosește tipuri obiectuale pentru descrierea datelor multimedia, astfel:

- tipul *ORDAudio* pentru datele audio,
- *ORDImage* pentru imagini statice,
- *ORDVideo* pentru secvențe video și
- *ORDDoc* pentru date eterogene.

Toate aceste tipuri conțin și informații despre sursa datelor, tipul de dată folosit fiind *ORDSource*.

Instanțele acestor tipuri conțin atribute, inclusiv metadate, metode și datele propriu-zise multimedia. Metadatele sunt folosite pentru stocarea informațiilor despre date, ca de exemplu mărimea fișierului multimedia, formatul de compresie, etc.

Definițiile tipurilor multimedia se găsesc în schema *ORDSYS*.

La inserarea datelor multimedia în baza de date Oracle sub controlul *interMedia*, indiferent de modalitatea de stocare a datelor multimedia, în interiorul sau exteriorul bazei de date, sunt automat extrase metadate pentru datele audio, imaginile statice și secvențele video. Metadatele extrase automat sunt:

- informații legate de stocarea datelor: locația și numele sursei datelor și locul unde sunt stocate datele, în baza de date sau sub forma unui fișier extern;
- data și ora curentă;
- formatul datelor;
- tipul în format MIME al datelor;
- pentru datele audio: formatul de codare, numărul canalelor audio, rata de eșantionare, tipul compresiei, durata secvenței;
- pentru imaginile statice: formatul imaginii, dimensiunile unui cadru, formatul de compresie al imaginii;

- pentru datele video: dimensiunile unui cadru, rezoluția unui cadru, rata de derulare a secvenței, durata secvenței, numărul total al cadrelor, formatul compresiei.

Funcționalitatea cheie într-o bază de date multimedia este modalitatea de regăsire eficientă a datelor multimedia continue și non-continue. O metodă general folosită pentru regăsirea bazată pe conținut a obiectelor multimedia se bazează pe simpla extracție a proprietăților obiectelor multimedia. În recunoașterea bazată pe conținut, se adaugă interpretarea semantică a obiectelor. Interpretarea semantică poate fi adăugată la indexare sau poate fi obținută printr-un proces de indexare semiautomat.

InterMedia suportă imagini digitale, statice, în două dimensiuni stocate ca și reprezentări binare ale unor obiecte sau scene ale lumii reale, utilizând cele mai populare formate de fișiere și scheme de compresie. Utilizatorii pot ușor să stocheze imagini în baza de date, fără să fie experți în caracteristicile diferitelor formate de fișiere pentru imagini, deoarece nu este necesar să se facă explicit conversii ale fișierelor la/de la un format intern folosit de Oracle pentru memorarea acestora. Oracle *interMedia* suportă următoarele tipuri de fișiere pentru imagini: TIFF, JFIF (JPG), BMP, TARGA, EXIF, PCX, PICT, GIF, CALS, SUN RASTER, FPIX, PNGF, PPMF, PGMF, PBMF, WBMP, RPIX. Utilizatorii sunt scutiți de cunoșterea complexității tehnologiilor de compresie/decompresie, ei trebuind doar să solicite ca imaginile să fie comprimate/decomprimate. Oracle *interMedia* suportă cele mai populare și eficiente scheme de compresie, incluzând CCITT G3/G4, ISO/CCITT JPEG, etc.

Recunoașterea imaginilor bazată pe conținut este o problemă importantă asociată sistemelor de gestiune a bazelor de date. Odată cu creșterea volumului colecțiilor imaginilor digitale care pot fi eventual stocate în baza de date, crește și dificultatea regăsirii imaginilor relevante. Pentru rezolvarea acestei probleme există două metode și ambele utilizează metadate pentru regăsirea imaginilor:

- folosind informații introduse manual în tabele, ca de exemplu: titluri, cuvinte cheie descriptive preluate dintr-un vocabular limitat și scheme de clasificare predefinite,
- folosind caracteristicile imaginilor, caracteristici extrase automat și recunoașterea obiectelor pentru clasificarea conținutului imaginii.

InterMedia permite combinarea celor două alternative prin proiectarea unei tabele ce conține imagini. Pentru aceasta se folosesc date de tip text pentru descrierea semnificației semantice a imaginii și tipul de dată *ORDImageSignature* pentru interogări bazate pe conținut ce folosesc atributele esențiale ale imaginii.

InterMedia este o extensie a SGBD Oracle introdusă începând cu versiunea 9i și oferă facilități pentru stocarea imaginilor, funcționalități de recunoaștere a imaginilor și facilități de conversie a formatului, prin introducerea unui nou tip de obiect *ORDImage* și a metodelor și funcțiilor asociate. Recunoașterea formelor este realizată prin posibilitatea de extragere a unui vector de caracteristici ale imaginii din diferite atribute vizuale.

Un sistem de recunoaștere bazat pe conținut prelucrează datele din imagine și creează o abstractizare a conținutului pentru atributele vizuale. Interogările lucrează cu abstractizarea imaginii și nu cu imaginea propriu-zisă.

Criteriile de căutare folosite de *InterMedia* sunt culoarea, textura și conturul și poziția. Pozițiile acestor atribute vizuale în cadrul imaginii sunt reprezentate prin coordonate. Aceste coordonate nu sunt folosite în mod independent pentru recunoașterea formelor ci doar împreună cu unul din cele trei atribute vizuale.

Imaginea odată inserată în baza de date este analizată și este stocată câte o reprezentare compactă a conținutului sub forma unui vector de caracteristici numit semnătura imaginii. Semnătura imaginii este extrasă prin segmentarea acesteia în regiuni, pe baza petelor de culoare care compun imaginea. Fiecare regiune are asociate informații despre culoare, textură și contur. Semnătura conține aceste informații pentru fiecare regiune care formează imaginea și informații despre culoare, textură și contur pentru reprezentarea acestor atribute în întreaga imagine.

Atributul *culoare* memorează informații despre distribuția culorilor în întreaga imagine. Această distribuție conține date despre intensitatea fiecărei culori.

Atributul *textură* reprezintă șabloanele din cadrul imaginii, precum granularitatea și netezimea. Spre deosebire de atributul contur, textura este foarte sensibilă la caracteristicile care apar cu mare frecvență în imagine.

Conturul este determinat de tehnicile bazate pe segmentare. Conturul este caracteristica unei regiuni de culoare uniformă.

Locația reprezintă poziția componentelor culoare, textură și contur. Recunoașterea imaginilor stocate într-o bază de date se face prin compararea lor cu o imagine model care poate fi o imagine stocată în baza de date, din afara bazei de date sau o imagine vectorială.

În procesul de căutare, se atribuie o *pondere* fiecărui atribut vizual în funcție de importanța lui. Valoarea fiecărei ponderi reflectă cât de sensibil trebuie să fie procesul de căutare față de un anumit atribut. Valorile ponderilor trebuie să fie între 0 (atribut nesemnificativ) și 1 (atribut extrem de important în procesul de căutare).

Asemănarea dintre două imagini pentru fiecare atribut vizual este calculată ca *scorul* sau *distanța* dintre imagini, respectând atributul. Scorul ia valori în intervalul 0 (nu există diferență) – 100 (diferență maxim posibilă). Scorul întregii imagini se calculează ca sumă a scorurilor atributelor ponderată cu importanța fiecărui atribut.

În procesul de căutare se folosește o valoare – *prag de semnificație*. Dacă suma ponderată este mai mică sau egală cu valoarea pragului atunci imaginile se potrivesc, altfel nu.

Pentru creșterea vitezei de căutare în bazele de date de mari dimensiuni care conțin date multimedia este utilă crearea și utilizarea unui index folosit pentru căutarea printre semnăturile imaginilor. Pentru aceasta se folosește un *index de domeniu* sau *index extensibil* deoarece acesta suportă obiecte complexe. Baza de date Oracle și *interMedia* cooperează pentru definirea, construirea și întreținerea unui index pentru datele de tip imagine, index de tip *ORDImageIndex*. Odată creat, indexul este automat actualizat ori de câte ori imaginile sunt inserate, modificate sau șterse din baza de date. Datele indexului sunt stocate în două *tablespace*-uri care trebuie create în prealabil: unul care conține datele indexului curent și celălalt este un index intern creat pe aceste date.

Recunoașterea formelor este o procedură complexă. Algoritmul de recunoaștere a formelor implementat în Oracle poate fi folosit numai în anumite condiții, și anume:

- dacă obiectul sau obiectele căutat(e) ocupă o parte semnificativă a imaginii;
- dacă nu există elemente irelevante suprapuse peste o parte a obiectului căutat;
- dacă obiectul căutat se află în aceeași parte a imaginii;
- dacă dimensiunile relative ale obiectului în cele două imagini, imaginea de referință și cea în care se face căutarea sunt apropiate;

- dacă obiectul căutat este fotografiat din același unghi în ambele imagini;
- dacă obiectele adiacente din imagine au culori distincte;
- dacă imaginea este formată doar din câteva forme simple.

Pentru a îndeplini aceste condiții, se pot decupa succesiv zone din imagine, zone în care se realizează căutarea și se pot utiliza diferite combinații de ponderi ale atributelor folosite în procesul de căutare.

Structura pentru indexare multidimensională *ORDImageIndex* permite accelerarea accesului la vectorul de caracteristici stocate.

De asemenea, *interMedia* poate gestiona date audio având următoarele formate: AIFF, AIFF-C, AUFF, WAV, MPEG I, MPEG II, MPEG III și formate audio Real Networks (prin intermediul serverului de streaming Real Networks) sau date video cu formatele: QuickTime, AVI, MPEG și formate video Real Networks (prin intermediul serverului de streaming Real Networks). El extrage în mod automat informațiile de tip metadata din aceste formate și le memorează în atributele obiectelor audio sau video din Oracle *interMedia*. Datele audio sau video propriu-zise pot fi stocate fie local în baza de date Oracle sau pot fi referite din surse externe menționate anterior. Sunt recunoscute schemele de compresie ADPCM și MU-LAW pentru datele audio și AVI Indeo pentru datele video.

Oracle *interMedia* combinat cu Oracle și cu unelte de dezvoltare a aplicațiilor furnizate de parteneri ai firmei Oracle, constituie o platformă puternică pentru dezvoltarea și întreținerea aplicațiilor Web de la nivelul întreprinderilor.

IBM DB2 Universal Database Extenders

DB2 UDB, serverul de baze de date al firmei IBM, dispune și el de tipuri de date noi care permit stocarea și gestionarea datelor multimedia. Pe lângă tipurile de date, așa numite clasice: numerice, șir de caractere, date calendaristice, etc., există și tipuri de date care permit stocarea și gestionarea datelor de tip audio, video, voce, figuri, etc. Acestea pot fi memorate în baza de date în câmpuri de date de tip BLOB (*binary-large-object*), șiruri binare mari de lungime variabilă care pot păstra date nestructurate cum sunt figurile, vocea, imaginile video, etc. De asemenea, există și tipuri de date care pot stoca date de tip șir de caractere de lungimi foarte mari numite CLOB (*character large object*

string) care pot stoca documente text, sau tipuri de date care permit stocarea unor șiruri grafice de lungime fixă sau variabilă care pot fi chiar de dimensiuni foarte mari precum GRAPHIC, VARGRAPHIC, LONG VARGRAPHIC, DBCLOB.

DB2 UDB este îmbogățit și cu componente numite *Extenders* pentru a permite gestionarea unor tipuri noi și complexe de date precum: imaginile, secvențele video, secvențele audio, obiectele spațiale, documentele text, etc.

IBM DB2 Universal Database Extenders extinde gestiunea datelor, incluzând suport pentru gestiunea imaginilor, a secvențelor video, audio și a obiectelor spațiale. Toate aceste tipuri de date sunt modelate, accesate și manipulate prin intermediul unui suport comun. Extensiile multimedia permit importul și exportul obiectelor multimedia și a atributelor acestora în interiorul și în afara bazei de date, controlând accesul la tipurile de date neconvenționale, cu același nivel de protecție ca în cazul datelor tradiționale și navigând prin sau extrăgând obiectele găsite din baza de date.

De exemplu, *DB2 Image Extender* definește un tip de dată distinct, *DB2IMAGE* ce are asociate funcții definite de utilizator pentru stocarea și manipularea fișierelor de tip imagine. Conținutul fișierului imagine pe care-l descrie *DB2Image* poate fi stocat ca *BLOB* sau în afara bazei de date, în sistemul de fișiere.

DB2 Image Extender oferă funcționalități de căutare similare bazate pe tehnologia QBIC (*Query by Image Content*) pentru imaginile stocate în tipul *DB2IMAGE*. Tehnologia QBIC oferă posibilitatea de interogare sau căutare a imaginilor pe baza conținutului lor.

DB2 NetSearch Extender oferă utilizatorilor și dezvoltatorilor de aplicații o metodă inteligentă de căutare a documentelor text care sunt stocate în baza de date sau în fișiere sistem externe. De asemenea, el permite gestiunea documentelor structurate de tip XML, HTML, GPP.

DB2 Spatial Extender permite memorarea și gestiunea datelor spațiale ale unor caracteristici geografice cum ar fi: obiecte geografice (de tipul râurilor, pădurilor, munților, etc.), o zonă bine definită în spațiu (de ex. o zonă de siguranță aflată în jurul unui obiectiv, o zonă în care operează un anumit business), un eveniment care apare la o locație ce poate fi definită (de ex. un accident auto care se produce într-o anumită intersecție). *DB2 Spatial Extenders* permite gestionarea acestor date spațiale asemănător celor tradiționale și folosirea lor în cadrul diferitelor tipuri de aplicații.

IBM Informix DataBlades

IBM Informix DataBlades este o extensie a serverului de baze de date Informix al firmei IBM care aduce puterea Internetului direct în baza de date. Astfel el permite:

- extinderea funcționalităților serverului dinamic IBM Informix;
- facilități specifice pentru gestiunea conținutului;
- facilități pentru customizare;
- facilități pentru gestiunea de fișe a conținutului multimedia.

Modulele lui *IBM Informix DataBlades* sunt extensii ale serverului care sunt integrate în interiorul motorului bazei de date, furnizând funcționalități noi pentru aplicații și crescând performanța sistemului. Cu ajutorul acestor module imaginile, documentele structurate, imaginile video complexe, paginile de informații HTML pot fi memorate, indexate și accesate în mod direct în baza de date, furnizând astfel un mai mare grad de performanță, securitate și ușurință în utilizare.

Există mai multe module ale *IBM Informix DataBlades*:

- *Spatial DataBlade* – modul care permite gestiunea inteligentă a informațiilor complexe, geospațiale, împreună cu date tradiționale, păstrând mecanismele eficiente de la nivelul modelului relațional al datelor;
- *Geodetic DataBlade* – modul care suportă interogări bazate pe date spațiale globale și pe date de timp fără limitările inerente în cazul proiecțiilor pe hărți și care asigură o precizie ridicată în legătură cu localizarea globală;
- *TimeSeries DataBlade* – modul care furnizează suport sofisticat pentru gestiunea datelor temporale și de tip serii-de-timp.
- *TimeSeries Real-Time Loader* – modul care este un loader de date care lucrează împreună cu modulele *TimeSeries DataBlade* și cu *NAG (Numerical Algorithms Group) DataBlade* pentru a atinge o performanță și o capacitate analitică superioară celei care este posibil să fie atinsă utilizând bazele de date relaționale tradiționale sau soft de sine stătător dedicat analizei de timp-real.

- *Web DataBlade* – modul care este o colecție de unelte, funcții și exemple care ușurează dezvoltarea aplicațiilor “inteligente” interactive de baze de date în medii Web.

Un punct important, dar dificil de realizat, este formularea și prelucrarea interogărilor complexe. Nici unul din sistemele menționate nu suportă căutări complexe, ca de exemplu imaginile care conțin persoane stând în fața unui autoturism de o anumită culoare. În acest caz, forma și culoarea autoturismului nu mai pot fi folosite pentru stabilirea asemănărilor dintre imagini deoarece în acest caz, o parte a autoturismului nu este vizibilă.

2.2. Comparație între formate

Trebuie luat în considerare faptul că scopul final este o prezentare a documentelor unor utilizatori interesați astfel că ieșirile din acest sistem ar trebui să fie documente ce să permită căutarea în textul acestora pe cât posibil fără a deteriora informația vizuală din documentele originale.

Acestea fiind spuse, devine destul de clar că documentele create trebuie să aibă următoarele caracteristici:

- trebuie să conțină atât imagini cât și text;
- să existe posibilitate de căutare;
- să fie citibile de către clienți web uzuali;
- să aibă conceptul de pagină;
- să fie de dimensiuni reduse;
- să fie suportate de către sistemele de recunoaștere - OCR existente (ca documente de ieșire).

Prin analiză criterială am stabilit importanța caracteristicilor și am luat în considerare mai multe tipuri de documente existente inclusiv un tip de document proprietar care ar putea fi dezvoltat special pentru acest proiect.

În urma analizei s-a ajuns la concluzia că dezvoltarea unui format proprietar nu este cea mai bună soluție, documentul PDF părând a fi cea mai bună soluție (tabel 4):

Tabel 4. Comparație între formate

Caracteristică	Importanța caracter.	TIFF	DOC	PDF	Format proprietar
Image + text	3	0	1	1	1
Cutare	4	0	0.5	1	1
Compatibilitate cu clienții existenți	5	0.75	1	1	0
Paginabile	2	1	0.75	1	1
Dimensiune	1	0.75	0.6	0.8	1
Suportabilitate	6	1	1	1	0
Punctaj acumulat		12.5	18.1	20.8	10

În plus față de celelalte documente, formatul PDF permite așezarea textului recunoscut sub imaginea inițială (text under the image) astfel că utilizatorul va vedea imaginile dar va căuta în text. Acest lucru asigură o calitate ridicată a vizualizării chiar și în cazul în care recunoașterea a fost de o calitate mai scăzută.

Documentele PDF permit și înmagazinarea de metainformații în acestea. Mecanismele de protecție ale documentelor PDF nu sunt de neglijat fiind posibil să protejăm un document la salvare, copiere, tiparire, copy-paste, etc.

Sistemul trebuie să fie gândit astfel ca toate documentele intermediare să fie păstrate, în cazul unor erori nemaifiind necesară reluarea întregului proces (de la scanare la salvare PDF) ci doar părți ale acestuia.

3. Conținutul web

3.1. Identificarea formatelor de documente pentru conținutul web

PDF - Portable Document Format

Portable document format PDF (Format portabil de documente) – este un format de fișier creat de firma Adobe în 1993 pentru schimbul de documente. PDF este utilizat pentru reprezentarea documentelor într-un spațiu bidimensional într-o manieră independentă de aplicațiile software, hardware și sistemele de operare.

Fiecare fișier PDF cuprinde o descriere completă despre un plan general pentru documentele cu două dimensiuni (sau cu trei dimensiuni) care cuprinde textul, fontul, imaginile și un vector grafic 2D cu conținutul documentului. Standardul PDF a fost oficial publicat doar în anul 2008.

Fundamente tehnice

PDF combina 3 tehnologii:

- un subset a limbajului de programare PostScript pentru descrierea paginii, pentru generarea graficelor și structurii de bază;
- un sistem pentru includerea / eliminarea fonturilor pentru a permite ca fonturile să existe în același document cu datele;
- un sistem structurat de memorare care să permită ca aceste trei elemente precum și orice conținut asociat să poată exista împreună într-un singur fișier, de asemenea să permită compresia datelor.

PostScript

Este un limbaj de descriere a paginii care rulează într-un interpretor pentru a genera imaginea. Acest proces necesită multe resurse. PDF este un format de fișier, nu un limbaj de programare (comenzile pentru controlul fluxului cum ar fi *if* și *loop* sunt eliminate, în timp ce comenzile grafice cum ar fi *lineto* sunt păstrate). Adesea PDF este generat dintr-o sursă PostScript. Comenzile grafice care sunt generate de codul PostScript sunt colectate și însemnate; orice fișier, grafic sau font la care se referă documentul sunt preluate și memorate de asemenea în fișier. Apoi totul este comprimat într-un singur

fișier. Așadar, toate informațiile generate de PostScript (font, structura de bază, dimensiunile) rămân nemodificate.

Ca format de document PDF-ul are câteva avantaje comparativ cu PostScript:

- PDF conține tokenizarea și interpretarea rezultatelor codului sursă PostScript, pentru o corespondență directă între modificările item-urilor din descrierea paginii PDF și modificările asupra modului de prezentare a paginii;
- PDF suportă transparență grafică în timp ce PostScript nu;
- PostScript este un limbaj de programare interpretiv cu o stare globală implicită, unele instrucțiuni care însoțesc descrierea unei pagini pot afecta modul de reprezentare a paginilor următoare. Așadar, toate paginile de procesat în documentul PostScript trebuie procesate în ordine pentru a determina modul de reprezentare corect pentru pagina dată, în timp ce fiecare pagină din documentul PDF este neafectată de caracteristicile celorlalte pagini. Ca rezultat, view-erele PDF permit utilizatorului să sară rapid la paginile finale dintr-un document lung, în timp ce un viewer PostScript necesită procesarea secvențială a tuturor paginilor înainte de a fi capabil să afișeze pagina cerută.

Structura fișierului

Fișierul PDF constă din primitive de *obiecte*, de 8 tipuri diferite:

- *Boolean* - valori booleene, reprezentând *adevărat* sau *fals*
- *Numbers* (numere)
- *Strings* (șiruri de caractere)
- *Names*
- *Arrays* – colecții ordonate de obiecte
- *Dicționare* – colecții de obiecte indexate prin Nume
- *Streams* – de obicei conțin mari cantități de date
- *NULL* – obiectul null

Obiectele pot fi fie *direct* (incluse într-un alt obiect) fie *indirect*. Obiectele indirecte sunt numărate printr-un *număr obiect* și printr-un *număr generat*. O tabelă de indecși numită *xref table* păstrează biții de ofset, pentru fiecare obiect indirect, față de

începutul fișierului. Această organizare permite un acces aleator eficient la obiectele din fișier și de asemenea permite ca mici modificări să poată fi efectuate fără rescrierea întregului fișier (actualizare incrementală). Începând cu versiune PDF 1.5, obiectele indirecte pot de asemenea să fie localizate într-un stream special cunoscut ca *object streams*. Această tehnică reduce dimensiunea fișierelor care au un număr mare de obiecte indirecte mici și este de ajutor pentru Tagged PDF.

Există două plane generale pentru fișierele PDF – neliniar (ne optimizat) și liniar (optimizat). Fișierele PDF neliniare consumă mai puțin spațiu pe disc decât partea de numărare liniară, așadar sunt mai lenți din punct de vedere a accesului deoarece porțiuni din datele necesare pentru a asambla pagina documentului sunt împrăștiate în tot documentul. Fișierele PDF liniare (de asemenea mai sunt numite „optimizate” sau „optimizare web”) sunt construite într-o manieră care asigură ca ele să fie citite într-un browser web (utilizând un plug-in), ele sunt scrise pe disc într-un mod liniar (cum ar fi în ordinea paginii). Fișierele PDF pot fi optimizate utilizând software Adobe Acrobat sau pdfopt care face parte din Ghostscript.

Modelul imaginii

Idea de bază despre reprezentarea graficelor în PDF este foarte asemănătoare cu cea din PostScript, excepție este doar utilizarea unui modul de transparență care a fost adăugată în PDF 1.4

Graficele PDF utilizează un sistem de coordonate carteziane independente de dispozitiv pentru a descrie suprafața unei pagini. O descriere a paginii PDF poate utiliza matrici pentru a scala sau roti elementele grafice. Conceptul cheie din PDF este acela că specificarea graficelor, care este o colecție de parametrii grafici care pot fi modificați, salvați și restaurați de descrierea paginii. PDF a avut 24 proprietăți grafice dintre care cele mai importante sunt:

- matricea transformării curente (CTM), care determină sistemul de coordonate;
- clipping path;
- spațiul culorilor;
- constanta alfa - care este o componentă cheie pentru transparență.

Graficele vectoriale

Graficele vectoriale în PDF, la fel ca în PostScript, sunt construite cu căi. Căile sunt de obicei compuse din linii și curbe cubice Bezier, dar pot fi de asemenea construite din liniile externe din text.

Similar cu PostScript, PDF-ul nu permite o singură cale pentru a grupa corect curbele și liniile. Căile pot fi o trăsătură, filtre sau utilizate pentru clipping. Trăsăturile și câmpurile pot utiliza orice set din grafic-ul selectat, incluzând *patter-nuri*.

PDF-ul suportă câteva tipuri de patternuri. Cel mai simplu este *paternul tiling* în care fiecare parte din șirul de descriere este specificată să desenată repetat. Aceasta poate fi un pattern de culoare tiling, cu culoarea specificată în obiect, sau un patern tiling necolorat în care diferă specificarea culorii în timpul în care tarenul este desenat.

Începând cu PDF 1.3 există de asemenea un patern *shading*, care desenează continuu variind culorile. Există 7 tipuri de patternuri *shading* dintre care cele mai simple sunt *radial shade*(Type 2) și *axial shade*(Type 3).

Imagini raster

Imaginile raster în PDF (numite imagini XObjects) sunt reprezentate prin dicționare cu stream-uri asociate. Dicționarul descrie proprietățile imaginii și stream-ul conține imaginea. (De obicei, o imagine raster poate fi inclusă direct în descrierea paginii ca și o imagine *inline*). Imaginile sunt de obicei *filtrate* în scopul comprimării. Filtrele de imagini suportate în PDF sunt următoarele:

- **ASCII85Decode** – un filtru vechi utilizat pentru a pune stream-urile pe 7-bit [ASCII](#);
- **ASCIIHexDecode** – similar cu cel de sus dar mai puțin compact;
- **FlateDecode** – un filtru uzual bazat pe [DEFLATE](#) sau algoritmul Zip;
- **LZWDecode** – un filtru vechi bazat pe compresia [LZW](#);
- **RunLengthDecode** – o metodă simplă de compresie pentru stream-uri cu date care se repetă utilizând algoritmul [Run-length encoding](#).

Filtrele specifice pentru imagini sunt:

- **DCTDecode** – un filtru cu pierderi bazat pe standardul [JPEG](#);

- **CCITTFaxDecode** – un filtru fără pierderi bazat pe metoda de compresie [CCITT fax](#);
- **JBIG2Decode** – un filtru cu sau fără pierderi bazat pe standardul [JBIG2](#) introdus în PDF 1.4;
- **JPXDecode** – un filtru cu sau fără pierderi bazat pe standardul [JPEG2000](#) introdus în PDF1.5;

În mod normal toate imaginile conținute în PDF sunt incluse în fișier. Dar PDF permite ca fișierele cu maginii să fie memorate în afara fișierului prin utilizarea *streamurilor externe* sau *imagini alternante*. Subseturile standardizate ale PDF, PDF/A și PDF/X interzic aceste tehnici.

Text

Textul în PDF este reprezentat prin *elemente text* în paginile care conțin streamuri. Un element text specifică ce *caracter* trebuie desenat la poziția respectivă. Caracterele sunt specificate utilizând *codificare* în funcție de fontul selectat. Un obiect font în PDF este descris ca un set de una sau mai multe litere, în funcție de dimensiunea fontului. Fișiere cu fonturi care pot fi incluse se bazează foarte mult pe formatul standard al fonturilor digitale: Type 1, TrueType și OpenType. În plus PDF suportă și Type 3 varianta în care componentele din fișier sunt descrise de operatorii grafici din PDF.

În timp ce caracterele text sunt afișate utilizând codurile caracterelor (numere întregi) pentru a mapa heroglifele în fontul curent se utilizează o codificare. Există un număr de codificări, incluzând *WinAnsi*, *MaRoman* și un număr mare de codificări pentru limbajele din Asia. Mecanismul de codificare din PDF este proiectat pentru fontul Type 1 și regulile utilizează reguli complexe pentru aplicarea lui în fonturile TrueType. Pentru fonturi mari sau fonturi cu heroglife nestandard, este utilizată codificarea specială *Identity-H* și *Identity-V*. În această codificare fiecare font trebuie să furnizeze o tabelă *ToUnicode* de informații semantice despre caracterele care vor fi prezente.

Transparența

Modelul de imagine original prezentat în PDF a fost similar cu cel din PostScript *opaque*: fiecare obiect este desenat pe pagină complet eliminând orice a fost desenat

înainte în aceeași locație. În PDF 1.4 modelul imaginii a fost extins pentru a permite transparența. Când este utilizată transparența, noile obiecte interacționează cu obiectele deja marcate pentru a crea mai multe efecte.

Extensia transparenței are la bază concepte cheie de *grupuri de transparență*, *modele amestecate*, *shap* și *alpha* modelul este apropiat de trăsăturile modelului Adobe Illustrator din versiunea 9. modelul amestecat are la bază modelul utilizat de AdobePhotoshop în acel moment. Când a fost publicată specificația PDF 1.4 formularele pentru calcularea modelelor amestecate au fost ținute secrete de Adobe. Fișierul PDF poate conține elemente interactive cum ar fi adnotațiile și câmpurile formei.

Securitate și semnătură

Fișierele PDF pot fi criptate pentru securitate, sau pot fi semnate digital pentru autentificare. Standardul de securitate furnizat de Adobe constă din două metode și două parole diferite „user password” și „owner password”. Documentul PDF poate fi protejat de parolă pentru deschidere (user password) și documentul poate fi de asemenea restricționat pentru anumite operații: listare, copiere text și grafic în afara documentului, modificarea documentului și adăugarea sau modificarea notelor text. Oricum, toate operațiile (excepție deschiderea documentului prin parolă de protecție, dacă există) pot fi restricționate prin parolă de „owner” sau „user”, sunt triviale datorită existenței softurilor „pdf crack”.

Variante de PDF

O mare parte din subseturile PDF au fost sau sunt standardizate pentru diferite componente:

- [PDF/X](#) – pentru imprimare și grafică ISO 15930
- [PDF/A](#) – pentru arhivare în medii ca întreprinderi / guvern]/ biblioteci]/etc ISO 19005
- [PDF/E](#) – pentru portabilitate și motorul de desenare (utilizabil din ISO TC171)
- [PDF/UA](#) – pentru fișiere PDF universal accesibile

HTM, HTML- Hypertext Markup Language

Unul din primele elemente fundamentale ale WWW (World Wide Web) este **HTML** (Hypertext Markup Language); acesta descrie formatul primar în care documentele sînt distribuite și văzute pe Web. Multe din trăsăturile lui, cum ar fi independența față de platformă, structurarea formătărilor și legăturile hipertext, fac din el un foarte bun format pentru documentele Internet și Web. Redus la esență, HTML (HiperText Markup Language - HTML) este un set de coduri speciale care se inserează într-un text, pentru a adăuga informații despre formatare și despre legături. HTML se bazează pe SGML (Standard Generalized Markup Language).

HTML a fost dezvoltat inițial de Tim Berners-Lee la CERN în 1989. Tim Berners-Lee a utilizat ca model SGML (un standard internațional în plină dezvoltare). SGML avea avantajul unei structurări avansate și al independenței de platformă dar proiectarea lui a avut în vedere mai mult structura semantică a documentului decît modul de formatare. Flexibil, SGML putea fi caracterizat ca o specificare pentru descrierea altor formate. Standardul oficial HTML este World Wide Web Consortium (W3C), care este afiliat la Internet Engineering Task Force (IETF). W3C a enunțat câteva versiuni ale specificației HTML, printre care și HTML 2.0, HTML 3.0, HTML 3.2, HTML 4.0 și, cel mai recent, HTML 4.01. În același timp, autorii de browsere, cum ar fi Netscape și Microsoft, au dezvoltat adesea în browserele lor. În unele cazuri, cum ar fi tagul Netscape, aceste extensii au devenit standarde *de facto* adoptate de autorii de browsere. Documentele HTML sînt documente în format ASCII și prin urmare pot fi create cu orice editor de texte. Au fost însă dezvoltate editoare specializate care permit editarea într-un fel de WYSIWYG deși nu se poate vorbi de WYSIWYG atîta vreme cît navigatoarele afișează același document oarecum diferit, în funcție de platforma pe care rulează. Au fost de asemenea dezvoltate convertitoare care permit formatarea HTML a documentelor generate (și formatare) cu alte editoare. Evident conversiile nu pot păstra decît parțial formătărilor anterioare deoarece limbajul HTML este încă incomplet.

Un fișier HTML este de fapt un fișier text. Este sursa paginii de Web. Acest fișier este scris în limbajul HTML sau HTM (ambele însemnând practic același lucru) pe care, la afișare, browserul îl interpretează. Orice fișier începe cu eticheta <html> și se termină cu eticheta </html>. Fișierul este compus din:

- **header**, inclus între etichetele <head> și </head>, conține informația introductivă de formatare a paginii
- **corp**, inclus între etichetele <body> și </body>

Structura generală a unui document HTML este următoarea:

<HTML>

<HEAD>

<TITLE> </TITLE>

</HEAD>

<BODY>

</BODY>

</HTML>

Între aceste etichete este inclus întregul document HTML. Ele comunică browserului unde începe și unde se termină documentul.

<HEAD> </HEAD>

Între aceste etichete sunt incluse titlul, deja menționat, precum și alte informații despre documentul HTML. Aceste elemente, numite *metatag*-uri sunt deosebit de importante pentru ca pagina să fie cât mai bine cotate de către motoarele de căutare. Ele vor face obiectul unui capitol separat.

Metatag-urile și antetul în ansamblul său nu sunt vizibile pentru vizitator în pagina Web. Totuși, la fel ca întreg codul HTML al paginii, antetul poate fi vizualizat selectând din meniul browserului opțiunea *View > Source*.

<TITLE> </TITLE>

Stabilește **titlul** documentului HTML.

Titlul documentului este deosebit de important deoarece este unul din criteriile folosite de motoarele de căutare pentru indexarea paginii.

<BODY> </BODY>

Conține descrierea textului și elementelor paginii. În corpul documentului se stabilesc, deci, aspectul și conținutul paginii Web. Elementele conținute în această secțiune sunt vizibile în pagină.

Limbajul HTML se bazează pe *tag-uri*. Acestea sunt comenzile, sau instrucțiunile pe care un browser le interpretează, ca să poată afișa corect pagina. Exemplu de tag: **<body>**. Semnele „<” (mai mic) și „>” se numesc *paranteze unghiulare*. Majoritatea tagurilor se folosesc în perechi; fiecare pereche are un tag de deschidere, (sau de început), și unul de închidere (de sfârșit). Tag-urile și textul dintre doua taguri se numește *element HTML*, iar textul dintre tag-ul de început și cel de sfârșit formează *conținutul elementului*. Exemplul de mai sus („<body>”) este tag-ul de deschidere pentru elementul „body”. Tag-urile de închidere au caracterul „/” în plus față de cele de deschidere. Scopul celor doua taguri „<body>” și „</body>” este să specifice browserului că elementul dintre ele este corpul paginii.

Sintaxa oricarui tag este:

<TAG atribut1="val" atribut2="val">Text</TAG>

Regulile de folosire a etichetelor, atributelor și valorilor acestora reprezintă **sintaxa** limbajului HTML. Similar limbajelor de programare, în HTML respectarea sintaxei limbajului este determinantă pentru o bună funcționare a documentului. Documentele HTML sunt interpretate de browser exact așa cum au fost ele scrise. Prin urmare, orice greșală de sintaxă se va reflecta direct în aspectul paginii Web, conducând, de cele mai multe ori, la o funcționare defectuoasă a acesteia. O pagina web poate conține:

- text
- imagini
- fișiere audio
- fișiere video

XHTML- Extensible Hypertext Markup Language

XHTML (eXtensible HyperText Markup Language) este un limbaj de marcare ce are aceleași capacități expresive ca și HTML, dar cu o sintaxă mai strictă. XHTML poate fi considerat ca încrucișarea dintre HTML și XML, în multe privințe fiind o reformulare a HTML în XML. La data de 26 ianuarie 2000, Consorțiul Web a făcut publică o nouă specificație-recomandare intitulată **XHTML**, reprezentând o familie, bazată pe XML, de tipuri de documente și module care extinde arhicunoscutul standard HTML. XHTML 1.0 este primul "descendent" al acestei familii, fiind în fapt reformularea tipurilor de documente HTML4 în termenii meta-lingajului XML. Dezvoltatorii de pagini și de aplicații Web vor migra astfel de la HTML 4, bazat pe relativ vechiul și complexul SGML (Standard Generalized Markup Language), la XML (eXtensible Markup Language) - extensibil și mai ușor de utilizat, cu un viitor întrevăzut a fi strălucit.

Arborele genealogic al XHTML

Avantajele unei aplicații XHTML sunt multiple, important este că pot fi citite de toate dispozitivele XML în timp ce păstrează compatibilitatea cu toate browserele de Internet mai vechi sau mai noi, fără a necesita specificații suplimentare.

Un document XHTML se compune din trei părți principale:

- declarația de conformitate sau DOCTYPE, definește tipul documentului creat. În cazul unui document XHTML, aceasta este `<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Strict//EN" "http://www.w3.org/TR/xhtml1/DTD/xhtml1-strict.dtd">`. Editorul HTML va ști să o introducă singur.
- secțiunea **<head>** conține meta-marcajele necesare clasificării și indexării corecte a paginilor pe internet. Urmează imediat după declarația de conformitate și este

delimitat în interiorul marcajului <head></head>. Cele mai importante meta-marcaje conținute în antet sunt <title></title> - titlul paginii, <meta name="description" content="" /> - descrierea paginii și <meta name="keywords" content="" /> - cuvintele-cheie care descriu cel mai bine conținutul paginii.

- secțiunea <body> , reprezintă pagina propriu-zisă, delimitată de marcajul <body></body>. Aici apar toate celelalte marcaje folosite.

Documentele XHTML pot fi etichetate atit "text/html" cât și "text/xml", elementul rădacină fiind <html>. Declarația tipului de documente XHTML se va face prin intermediul construcției DOCTYPE, ca în SGML, existând trei tipuri de definiții de documente conform specificației HTML 4: tipul *strict* , *tranzitional* și *pentru cadre (frames)*, ca în exemplul următor:

```
<!DOCTYPE html
PUBLIC "-//W3C//DTD XHTML 1.0 Strict//EN"
"DTD/xhtml1-strict.dtd">
```

Structura de baza a unui document este următoarea (fig.1):

```
<!DOCTYPE ...>
<html>
<head>
<title>... </title>
</head>
<body> ... </body>
</html>
```

XHTML folosește *tag-uri* pentru a structura documentul. Deosebirile față de HTML (deosebiri care derivă din faptul că XHTML este bazat pe XML) sunt următoarele:

- Tag-urile (și atributele lor) se scriu cu litere mici (XML este case-sensitive)
- Toate tag-urile trebuie închise
- Valorile atributelor trebuie închise în ghilimele

XHTML este un limbaj care folosește markere și are aceeași putere de expresie ca și HTML dar mai are în plus un avantaj, acela de a respecta sintaxa XML. Pe când HTML este o aplicație de tip SGML (Standard Generalized Markup Language), care este un

limbaj foarte flexibil, XHTML este o aplicație a XML care este un subset mult mai restrictiv al SGML. Documentele XHTML trebuie să fie bine formate, de aceea ele permit ca procesarea automată să fie făcută cu unelte standard XML, spre deosebire de HTML care necesită un parser destul de complex, tolerant și personalizat.

Scopul folosirii XHTML-ului este acela de a înlocui limbajul clasic HTML, de a schimba tehnologia folosită cu una mai performantă, realizând un pas important în dezvoltarea web-ului; în plus, este susținut de către organizația World Wide Web Consortium (W3C).

XML: Extensible Markup Language

XML (eXtensible Markup Language), descendent al SGML (Standard Generalized Markup Language), este un meta-limbaj utilizat în activitatea de marcare structurală a documentelor, a cărei specificație a fost dezvoltată începând cu 1996 în cadrul Consorțiului World Wide Web (W3C), de un grup de cercetare condus de Jon Bosak de la Sun Microsystems, la care au aderat ulterior o serie de grupuri de experți din comunitățile academice (Text Encoding Initiative, NCSA, James Clark) și industriale (SUN, IBM, Netscape, Oracle, Adobe ,etc.).

Acest limbaj nu procesează în nici un fel datele sau informația. Menirea lui principală este doar de a forma și structura informația. Principalul avantaj al XML este compatibilitatea. Informația structurată cu ajutorul XML va fi citită și interpretată ulterior în același fel, indiferent de dispozitiv, fie el mobil, palmtop, PC sau Mac. XML a fost elaborat pentru:

- separarea **sintaxei** de **semantică** pentru a furniza un cadru comun de structurare a informației;
- **construirea de limbaje de mark-up** pentru aplicații din orice domeniu;
- **structurarea informației** în viitor;
- **asigurarea independenței** de platformă și suport pentru **internaționalizare**.

XML este un format de fișier asemănător cu Microsoft Word sau un fișier Adobe Acrobat, o foaie de calcul tabelar sau o pagină online HTML. Oricum, are proprietăți diferite de acestea și anume:

- standardul XML nu este controlat de o singura organizație; este un 'open' standard la care oricine poate contribui sau folosi;
- fișierle XML nu sunt salvate ca date 'binare', ci mai degrabă ca plain text. Asta înseamnă că sunt independente de platformă și pot fi citite de oameni;
- fișierele XML nu au o structură specifică, însă folosesc un set de reguli de bază (dar stricte). Astfel, se pot reprezenta prin fișiere XML multe tipuri de date și informații, de la documente până la fișiere cu imagini și tranzacții financiare;
- regulile XML pot fi folosite să restricționeze structura tipurilor de date - deci pot fi create noi standarde. Structura este de tipul self- descriptive și fiecărei date îi este asociat un 'tag' pentru a o descrie într-un fel. Acest lucru permite ca fișierele să poată fi validate de un calculator, dar chiar și de oameni până la un anumit nivel; sunt doar secvențe de text create pentru a structura, depozita și transporta informație. Important este că definind reguli stricte dar neimpunând restricții asupra structurii, se pot dezvolta formate 'standard XML' (numite Schemas) care să reprezinte un tip particular de date. Există "Standard XML schemas" pentru multe tipuri de informații, de la tranzacții de vânzări în business, până la formatarea știrilor.

Documentele XML sunt realizate din unități de stocare numite entități, ce conțin date parsate sau neparsate. Datele parsate sunt realizate din caractere, unele dintre ele formând date caracter iar altele ca marcaje. Marcajele codifică o descriere a schemei de stocare a documentului și structura logică. XML furnizează un mecanism pentru a impune constrângeri asupra schemei de stocare și a structurii logice.

Cuvantul *marcaj (markup)* a fost folosit inițial pentru a descrie anumite adnotări, note marginale în cadrul unui text cu intenția de a indica tehnoredactorului cum trebuie listat un anumit pasaj. Generalizând, putem defini marcajul drept orice acțiune de a interpreta explicit o porțiune de text.

Un *marcaj (tag)* este un șir de caractere delimitat de caracterele "<" și ">". *Datele caracter* reprezintă conținutul marcajelor.

Coduri fișier ". XML"

Un document XML este un **arbore ordonat etichetat**:

- **date caracter** - noduri frunză ce conțin datele
- noduri **elemente** etichetate cu
 - un nume (adesea numit și tipul elementului)
 - o mulțime de atribute, fiecare din ele având un nume și o valoare

XML-ul este un limbaj cu o sintaxă simplă și care permite doar structurarea datelor într-o manieră proprie prin definirea propriilor taguri. Această facilitate de structurare a datelor a permis folosirea sa pentru a dezvolta limbaje noi precum și pentru a fi folosit în noi standarde de stocare a datelor. Nu în ultimul rând, XML-ul poate fi folosit pentru a schimba date între aplicațiile care au nevoie de a comunica într-un limbaj comun.

ASP- Application Server Provider

Există câteva definiții date termenului de **ASP (Application Service Provider** - furnizor de aplicații livrate ca servicii), printre care și următoarele două, date organizații de specialitate:

- **Asociația Americană pentru IT** oferă următoarea definiție:

"Un ASP este orice companie care distribuie și întreține automat aplicații și servicii web pentru clienți/utilizatori, direct prin Internet sau printr-o rețea privată."

- **Aberdeen**, o binecunoscută companie de analiză de marketing și poziționare pe piața IT, definește ASP-ul astfel:

"ASP-ul este un tip de afacere care oferă consumatorilor clienților acces la anumite software-uri găzduite într-o locație administrată centralizat."

- Conform dicționarului online **Webopedia** ASPurile sunt:

"...terțe părți care întrețin și distribuie clienților soluții și servicii de tip software dintr-o locație administrată centralizat."

Toate aceste definiții pun accent pe faptul că ASP-ul **distribuie și întreține aplicații software prin Internet**. Cu alte cuvinte, un ASP "închiriază" aplicații software pe Internet clienților care nu pot sau nu doresc să investească în licențe pentru software, hardware, personal IT specializat și alte resurse, oferindu-le acestora posibilitatea de a prelua controlul total asupra software-ului. Conform sitului www.ASPnews.com, ASPurile se împart în 4 subcategorii:

- ASPuri pentru corporații - oferă aplicații high-end pentru afaceri;
- ASPuri locale/regionale - oferă numeroase tipuri de aplicații pentru afaceri mai mici (locale sau regionale);
- ASPuri specializate - oferă aplicații pentru anumite nevoi, de genul serviciilor pentru website-uri sau de resurse umane;
- ASPuri pentru piețe verticale - programe utilitare pentru un anumit tip de industrie, cum ar fi cea imobiliară.

Un ASP trebuie să combine cunoștințele legate atât de servicii, networking, tehnologie și management, cât și de administrarea, integrarea, managementul și suportul aplicațiilor. Din punct de vedere al utilizatorului, el cumpără un singur serviciu de la un singur distribuitor: ASP-ul. În realitate, există un sistem complex care susține ASP-ul și care oferă diferite calificări în diferite domenii de tipul networking-ului sau al aplicațiilor. În marea majoritate a cazurilor, componentele și serviciile ASP-urilor sunt distribuite printr-o serie de parteneriate între diferiți furnizori. Cel mai simplu mod de distribuție poate include un singur creator de programe, susținând propria aplicație, pe care o distribuie pe Internet sub termenii unui contract.

Cel mai la îndemână exemplu de serviciu ASP este webmailul. Chiar dacă mailurile pot fi accesate și printr-un program specializat (ex: Outlook, Firebird, etc.), marea majoritate a utilizatorilor persoane fizice preferă să se folosească de serviciul web-based (Yahoo, Gmail, Hotmail, etc.). Avantajele unui astfel de serviciu: posibilitatea de a-l accesa de pe orice computer, lipsa unor cerințe software speciale, faptul că nu necesită mentenanță. Aceste avantaje se mențin în cazul oricărui tip de serviciu ASP. Pentru acele

servicii profesionale oferite contra cost, un alt mare avantaj este partajarea costurilor de dezvoltare între toți utilizatorii sistemului. Serviciile web distribuite în regim ASP (application service provider) reprezintă una dintre laturile cele mai avansate ale Internetului. În esență, o astfel de aplicație oferă servicii online cu nivel de complexitate ridicat, servicii care se adresează unui target larg (ex: utilizatorii de mail, persoanele tinere, companii mici și mijlocii indiferent de domeniu) sau restrâns (avocați, agenții imobiliare, agenții turistice, contabili, etc.).

ASP-urile folosesc ca platformă de distribuție a serviciilor un sistem bazat pe soluții web. Fig. 2 arată cum utilizatori aflați online se pot conecta prin Internet la serverele web ale ASP-urilor localizate în spatele unui firewall. Serverul de web accesează informațiile din baza de date printr-o rețea privată.

ASP: Active Server Pages

ASP (Active Server Pages) este una dintre metodele de construire dinamică a paginilor de Web. Începuturile WWW-ului au constat în livrarea de pagini HTML prin protocolul HTTP. Acestea constau în text formatat, imagini și link-uri care permiteau navigarea de la o pagină la alta. În curând s-a simțit nevoia interacțiunii, posibilității de a prelua răspunsuri de la client, de a furniza pagini personalizate, de a afișa date prelucrate din diverse surse care nu sunt statice. Astfel, a apărut partea dinamică a Web-ului, parte ce a dus la o dezvoltare substanțială a Internetului prin e-commerce, aplicații online de tot felul, de la simpla completare a unor formulare electronice până la site-uri care simulează burse, bănci. În prezent, prin pagini dinamice de web se pot crea aplicații complexe care înlocuiesc din ce în ce mai mult aplicațiile clasice client-server pentru că sunt simplu de întreținut, nu necesită instalare pe calculatoarele client și pot fi accesate de oriunde. În acest cadru, ASP este una din opțiunile principale pentru crearea de pagini dinamice. ASP-urile se utilizează pentru:

- Editarea dinamică, schimbarea sau adăugarea de conținut la o pagină Web;
- Generarea de răspunsuri cererilor utilizatorilor sau informațiilor trimise prin formularele HTML;

- Accesarea oricărei date sau baze de date și returnarea de rezultate către navigator;
- Personalizarea paginilor Web, făcându-le mai utile pentru utilizatorii particulari;
- Simplitatea și viteza de execuție a ASP față de CGI și Perl;
- Asigurarea securității, deoarece codul ASP nu este vizibil prin navigator;
- Posibilitatea de a fi afișate de orice navigator, fișierele ASP sunt returnate în format HTML;
- Minimizarea traficului prin rețea.

La dezvoltarea unui site care utilizează ASP-uri pot interacționa aproape toate tehnologiile Web dinamice existente. În fig.1 este prezentat modul în care se integrează diferitele componente de sistem cu aplicațiile.

Caracteristica majoră a ASP constă în capacitatea sa de includere a scripturilor direct într-un fișier pe care-l accesează navigatorul, generând astfel *pagini dinamice*.

Dar, față de celelalte scripturi de pe server, fișierele ASP pot să conțină cod HTML sau XML, incluzând chiar scripturi pentru client și referiri la componentele COM (Component Object Model) prin care realizează o serie de sarcini, cum ar fi conectarea la o bază de date.

- text;
- tag-uri HTML sau XML;
- comenzi ASP.

CGI-Computer generated imagery

Computer-generated imagery (CGI) **CGI (Imagini generate pe computer)** este aplicarea domeniului graficii pe calculator (mai exact 3D Computer Graphics) și a efectelor speciale în filme, programe de televiziune, reclame, jocuri video. CGI a fost folosit prima dată în filme, în anul 1973. Progresele în sfera CGY sunt raportate în fiecare

an la Siggraph, o conferință anuală pe tema graficii pe calculator si tehnicilor interactive, care adună la Los Angeles mii de profesioniști din domeniu.

CGI: Common Gateway Interface

CGI (Common Gateway Interface) este un protocol standard de comunicare între documentele Web și aplicațiile localizate pe serverul Web. Este metoda prin care un server web poate să prezinte vizitatorilor informația prin intermediul unui browser web, în funcție de cererile acestora, interactiv, accesând baze de date sau documentație, și extrăgând datele cerute.

CGI (Common Gateway Interfaces) s-a impus ca cea mai eficientă, stabilă și ușor de înțeles modalitate de manipulare a informației generate în mod dinamic pe Web. Este de fapt acea parte a server-ului Web care poate comunica cu alte programe care rulează pe sistem. Cu ajutorul acestei interfețe, serverul Web poate apela un program. Comunicarea prin Internet se bazează pe protocoale. Protocolul utilizat pentru serverele Web, pentru comunicațiile dintre servere și programele care le accesează (clienți) este **HTTP (HyperText Transfer Protocol)**. Ca orice protocol pe Internet, el este în același timp un mecanism extrem de complex și un dialog simplu bazat pe principiul cerere-răspuns: programul care cere (clientul) emite o cerere spre server (după o secvență inițială de conectare). Cererea pe WWW pornește în general începând de la un document în format HTML și poate să ia o mare varietate de forme: de la o cerere simplă, de exemplu aducerea unui alt document, până la completarea form-urilor HTML. Datele transferate de către protocol pot veni în formate diferite, de la text simplu până la imagine sau videoclip. Mai mult, este posibil ca datele ce compun un document să se afle fizic pe calculatoare diferite, putând fi accesate doar prin alte protocoale de comunicație, nu direct via HTTP. Comunicația se derulează între clientul care a inițiat cererea și serverul Web accesat. Serverul are la dispoziție mecanismul specificat de standardul CGI. CGI specifică exact sub ce formă trebuie date mai departe informațiile sosite sub forma unei cereri de la client. Programele care recepționează cererile se numesc programe CGI: sunt practic programe activate pe server, indirect, de către utilizatori externi, datorită caracterului interactiv (dinamic) al documentelor HTML. Aceste programe pot fi scrise în principiu în orice limbaj de programare de uz general. CGI specifică modul în care

trebuie să se facă transferul cererilor sosite la serverul Web(prin HTTP) către programele CGI. CGI este deci o interfață independentă de limbaj, ce permite realizatorilor aplicațiilor Web să genereze *documente dinamice*. Numite uneori impropriu *scripturi CGI*, acestea pot fi scrise în aproape orice limbaj ce are posibilități de acces la variabilele de mediu și poate produce o ieșire.

Scripturile CGI trebuie să aibă o anumită extensie, de obicei extensia implicită este .cg dar poate fi și .pl, .js , sau oricare alta extensie, în funcție de limbajul de programare folosit la realizarea acestor scripturi. Pentru ca un programator HTML să poată folosi programe cgi acesta trebuie să ceară drepturi de acces administratorului serverului Web la directorul cgi-bin. Așa cum am mai specificat, programele CGI pot fi scrise în orice limbaj de programare care folosește intrări și ieșiri standard (STDIN, STDOUT). Limbajele cu ajutorul cărora se pot realiza programe CGI se împart în două categorii, și anume: programe compilate C, C++, Turbo Pascal, Fortran, Ada, etc. sau programe interpretate PERL, AppleScript, shell-urile UNIX, etc.

Scopul CGI este furnizarea unui mecanism flexibil și convenabil pentru extinderea funcțiilor unui server HTTP peste limitele unor simple preluări și afișări de fișiere.

Cele mai răspândite aplicații ale CGI sunt:

- prelucrarea datelor inserate în formulare (care necesită un răspuns);
- interogarea unor baze de date pentru o anumită informație (se realizează cu motoarele de căutare: Altavista, Lycos, Yahoo, Infoseek prin interogări SQL);
- documente virtuale (documente HTML complexe, care conțin text, imagini, fișiere de sunet sau video).

PHP- Hypertext Preprocessor

PHP (Hypertext Preprocessor) este unul din cele mai folosite limbaje de programare server-side. Limbajul PHP s-a născut în 1994 din nevoia lui Rasmus Lerdorf de a afla câte persoane îi vizitează CV-ul online. El a denumit setul de scripturi create PHP, acronimul pentru Personal home page. Pe parcursul următorilor trei ani limbajul a evoluat, adevăratul succes cunoscându-l când Zeev Suraski și Andi Gutmans au rescris motorul PHP de la cap la coadă, motor care poartă din versiunea 4 a PHP numele Zend, o

combinație de litere din prenumele creatorilor săi: Zeev si Andi. Fiind open-source, PHP beneficiază de suport activ din partea comunității online, acesta fiind și motivul creșterii explozive a numărului site-urilor bazate pe PHP.

Se folosește în principal înglobat în codul HTML dar începând de la versiunea 4.3.0 se poate folosi și în mod „linie de comandă” (CLI), permițând crearea de aplicații independente. Este unul din cele mai importante limbaje de programare web open source și server-side existând versiuni disponibile pentru majoritatea web serverelor și pentru toate sistemele de operare. Conform statisticilor este instalat pe 20 de milioane de situri web și pe 1 milion de servere web.

Caracteristici:

- **Familiaritatea** : sintaxa limbajului este foarte ușoară combinând sintaxele unora din cele mai populare limbaje Perl sau C;

- **Simplitatea** : sintaxa limbajului este destul de liberă. Nu este nevoie de includere de biblioteci sau de directive de compilare, codul PHP inclus într-un document executându-se între marcajele speciale;

- **Eficiența** : PHP-ul se folosește de mecanisme de alocare a resurselor, foarte necesare unui mediu multiutilizator, așa cum este web-ul;

- **Securitate** : PHP-ul pune la dispoziția programatorului un set flexibil și eficient de măsuri de siguranță;

- **Flexibilitate** : fiind apărut din necesitatea dezvoltării web-ului, PHP a fost modularizat pentru a ține pasul cu dezvoltarea diferitelor tehnologii. Nefiind legat de un anumit server web, PHP-ul a fost integrat pentru numeroasele servere web existente;

- **Gratuitate** : este probabil cea mai importantă caracteristică a PHP-ului. Dezvoltarea PHP-ului sub licența open source a determinat adaptarea rapidă a PHP-ului la nevoile web-ului, eficientizarea și securizarea codului.

PHP folosește extensii specifice pentru fișierele sale: .php, .php3, .ph3, .php4, .inc, .phtml. Aceste fișiere sunt interpretate de către serverul web, iar rezultatul este trimis în formă de text sau cod HTML către browser-ul clientului. Când accesăm o pagină HTML serverul care o gazduiește trimite pagina HTML către browser spre afișare. În cazul unei pagini PHP serverul citește codul PHP, îl interpretează și generează dinamic pagina HTML care este trimisă browserului spre afișare. Acesta este motivul pentru care utilizatorii folosesc PHP pentru construirea unor pagini cu conținut dinamic.

Când PHP-ul parcurge un fișier, de fapt "citește" textul până când întâlnește una din etichetele speciale care îi spun să înceapă să interpreteze textul ca pe cod PHP. Se execută codul până când este întâlnită eticheta de închidere. Apoi se "citește" din nou textul mai departe. Acesta este motivul pentru care se poate adăuga cod PHP în interiorul HTML-ului. Codul PHP este delimitat de unul din următoarele seturi de etichete de deschidere și închidere:

<?php	?>	etichete recomandate
<script language="php"?>	</script>	
<?>	?>	folosirea lor necesită anumite setări pe server
<%>	%>	etichete tip ASP, folosirea lor necesită anumite setări pe server

RSS: Formate de fluxuri web

RSS este o familie de fluxuri web, specificate în XML, și folosite pentru a distribui știri, noutăți sau sumaruri de către portalurile care oferă asemenea servicii.

Un flux web (sau „feed”) este un format de date, folosit pentru a oferi utilizatorilor conținut (în formă integrală sau în rezumat) împreună cu link-ul către sursă și o serie de elemente descriptive. Se mai numește și flux RSS, feed RSS, webfee, stream RSS sau canal RSS. Fluxurile web oferă acest conținut cel mai adesea sub forma unui fișier XML.

Apărut în 1997 și utilizat pentru prima dată de Netscape, feed-ul RSS a devenit un mod foarte popular de distribuie a informațiilor. Practic, cam toate site-urile cu adevărat importante ce oferă material updatat zilnic folosesc serviciul RSS.

Fluxul de informații de tip RSS este marcat pe site-urile ce oferă acest serviciu într-un mod ușor de identiifcat. Abonarea la fluxurile de tip RSS implică fie existența unei aplicații desktop destinată acestui lucru (feed reader, de genul client email - MS Outlook, Mozilla Thunderbird, Apple Mail - sau browser web - Internet Explorer, Mozilla Firefox,Apple Safari), fie o aplicație sau site web (Google Reader, MyYahoo, Bloglines, etc).

Modalitatea de abonare la un flux RSS al unui site depinde foarte mult de clientul de RSS folosit; de pildă, MS Outlook 2007 sau Mozilla Firefox:

Abrevierera RSS provine din referirea la următoarele standarde:

- * **Really Simple Syndication;**
- * **Rich Site Summary;**
- * **RDF Site Summary.**

RSS este un fișier XML; el descrie conținutul unui site, fiind actualizat odată cu acesta. Inițial, utilizat de Netcape pentru a crea pagini *My Netscape*, RSS a fost adoptat de serviciile de informare (news syndication services), weblog-uri. RSS este important pentru că permite informarea despre actualizări ale resurselor online fara a face apel la e-mailuri sau Newslettere, implicit reducând posibilitatea de a primi spam și viruși. Un fișier RSS include un logo, linkul la site, un câmp de citire și itemii de noutate; fiecare item constă din URL-ul articolului respectiv, titlul, o descriere . Sisteme numite agregatoare - *aggregators* sau *harvesters*, citesc RSS-urile și informează abonații lor asupra modificărilor de pe site, informația fiind astfel transmisă unei audiențe largi - *syndicated*. Informația din fișierele RSS este stocată in baze de date, care pot fi accesate cu multiple criterii de căutare.

Utilizatorii care vor să afle noutățile apărute pe site vor folosi un cititor de RSS ("RSS reader"), cititor de feed ("feed reader") sau un agregator ("aggregator"), care

poate fi instalat pe calculator sau poate fi integrat în browser (precum Feedster sau Blogdigger). Există agregatoare care permit gruparea mai multor surse de știri într-una singură. Utilizatorul se înscrie la un RSS prin introducerea URL-ului și astfel este început procesul de subscriere. Cititorul de RSS verifică în mod regulat feed-urile utilizatorului, downloadează update-urile pe care le gasește și oferă o interfață pentru a citi feed-urile.

SVG: Scalable vector graphics

Lista de programe care poate deschide fișiere cu extensia *.SVG*:

- firefox.exe;
- gimp-2.6.exe;
- ImgBurn.exe;
- Inkscape.exe;
- Paint Shop Pro 9.exe;
- sbrowser.exe.

SVG (Scalable Vector Graphics) este un limbaj pentru descrierea de imagini 2D folosind XML (Limbajul XML - *eXtensible Markup Language* - este limbajul ce oferă un format pentru stocarea și transmiterea de date). Acest format de fișier este folosit pentru salvarea de fișiere vectoriale grafice. Fișierele SVG pot fi vizualizate pe web și sunt susținute de browsere moderne care permit ca imaginile vizualizate să aibă o înaltă calitate grafică.

Este un standard al organizației W3C a cărui proiectare a început în anul 1999. Permite definirea imaginilor prin 3 metode: text, grafică vectorială și bitmap-uri.

Deși există aplicații specializate pentru crearea și editarea de SVG-uri, în acest scop poate fi folosit orice editor text. Vizualizarea unei imagini *SVG* poate fi realizată cu orice browser modern.

În momentul de față *SVG*-ul are mai multe profile pentru a se adapta mai bine la diferite constrângeri. Astfel, profilele SVG Tiny și SVG Basic au fost create special pentru dispozitivele mobile cu resurse limitate, iar profilul SVG Print este destinat mediilor de printare a documentelor. Pentru animarea unei imagini SVG organizația W3C recomandă standardul SMIL. Pe lângă recomandarea oficială mai există și alte soluții, precum ar fi ECMAScript.

Principalele elemente din componența unui fișier *SVG* sunt:

- Paths – „Căi” ce pot fi folosite pentru descrierea conturului unei forme. Conturul poate rămâne gol sau poate fi umplut. Căile pot fi folosite și pentru a specifica zonele de decupare;

- Forme de bază – Specificația *SVG* oferă posibilitatea folosirii formelor de bază: dreptunghi, cerc, elipsă, linii și poligoane. Acestea pot fi construite și cu ajutorul căilor care au același contur;

- Text – Pentru a specifica textul ce apare într-o imagine trebuie folosite elemente de tipul text;

- Painting – Se referă la posibilitatea de a umple formele specificate în *SVG*. Pentru aceasta poate fi folosită o singură culoare, o culoare cu transparență, un gradient sau un model;

- Culoare - Proprietatea culoare este folosită pentru specificare culorii;

- Gradient și model – Folosite pentru colorarea formelor specificate;

- Decupare, mascare - Pot fi folosite zone de decupare sau de mascare;

- Filtre – Descriu diferite efecte aplicate imaginilor;

- Interactivitate – O imagine *SVG* are posibilitatea de a interacționa cu utilizatorul. Astfel, la apăsarea unui buton sau la folosirea mouse-ului pot fi pornite diferite scripturi;

- Linkuri – Un document poate conține legături către alte pagini sau elemente web;

- Scripting – Într-un *SVG* pot fi definite scripturi cu diverse funcții;

- Animații – Pentru un *SVG* pot fi specificate diverse tipuri de animații;

- Font - Nu este necesar ca utilizatorul final să aibă instalate diferitele seturi de caractere folosite; fonturile pot fi incluse în imagine;

- Metadata – Pentru integrarea mai bună se oferă și opțiunea specificării de metadata (acestea sunt datele care descriu datele propriu-zise).

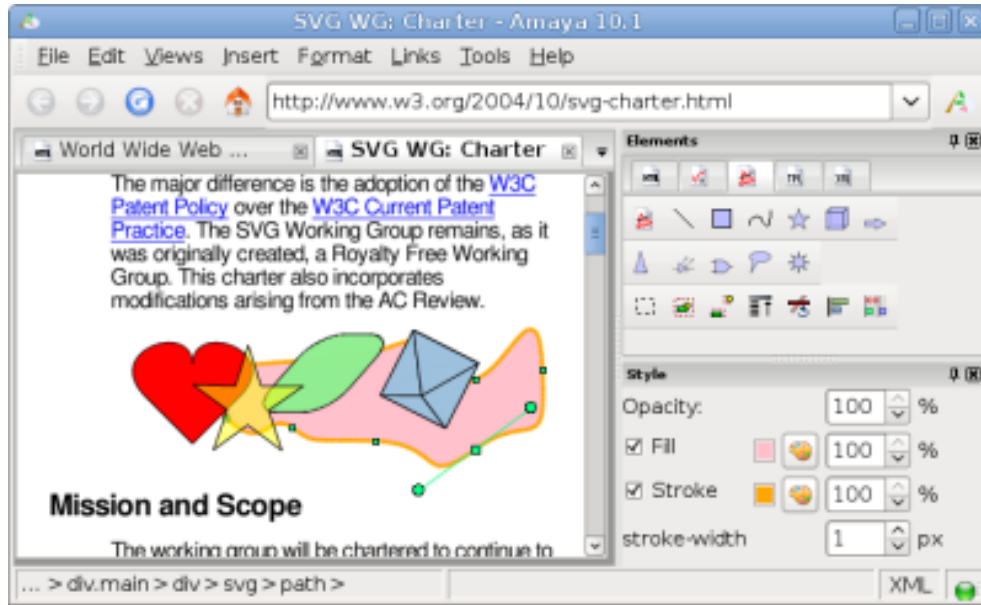


Fig.3.1. Formatul SVG

SVG este util pentru prezentare și comunicare, reprezentarea datelor statistice, științifice și medicale, pentru acces la informație pentru persoane cu dizabilități etc.

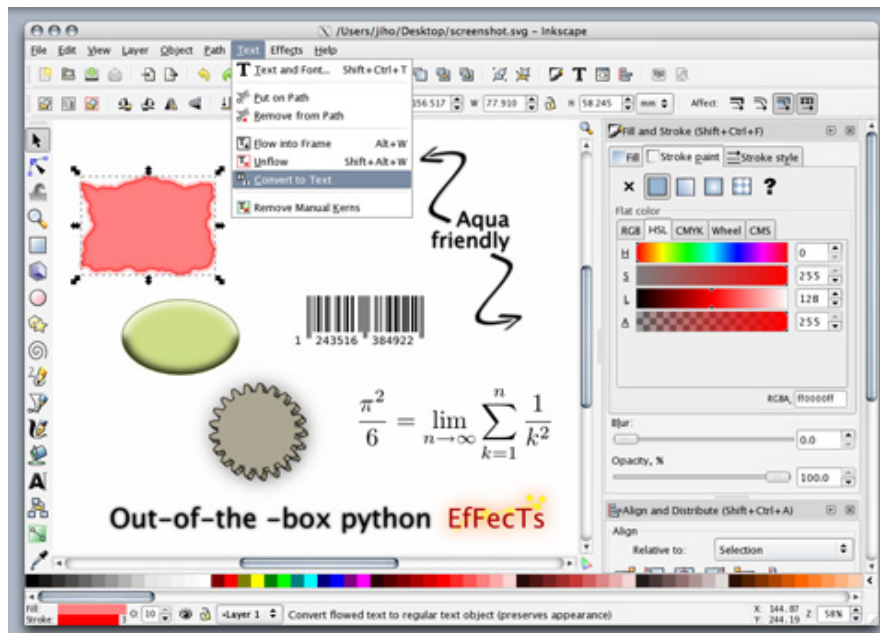


Fig.3.2. Formatul SVG grafice, formule

VRML: Virtual Reality Modeling Language

VRML (Virtual Reality Modeling Language)

...este un format standard de fișier care prezintă imagini 3D a diferitelor forme într-o mediu interactiv.

Fișierele *VRML* sunt de obicei numite "lumi" și au extensia WRL (de exemplu, island.wrl - <http://yallara.cs.rmit.edu.au/~rotaylor/web3d/journal.shtml>).

Prima versiune a *VRML* a fost specificată în noiembrie 1994. Versiunea actuală și funcțională este *VRML97* (ISO / IEC 14772-1:1997); *VRML* a fost acum înlocuit de *X3D* (ISO / IEC 19775-1). Termenul *VRML* a fost inventat de Dave Raggett într-un document depus la Prima Conferință Internațională de World-Wide-Web în anul 1994, dar el a fost dezvoltat de Tony Parisi și Peter Kennard.

În 1997 a fost finalizată o nouă versiune a formatului și *VRML97* (cunoscut ca *VRML2* sau *VRML 2.0*) și a devenit un standard [ISO](#). *VRML97* a fost folosit pe internet pe unele homepages personale și site-uri, cum ar fi "CyberTown", care a oferit chat 3D utilizând Blaxxun Software-ul. Formatul a fost promovat de către Cosmo SGI Software-ului; atunci când SGI a fost restructurat în 1998, divizia a fost apoi preluată de Computer Associates care a decis să nu dezvolte sau să distribuie software-ul.

Consortiul *VRML* a schimbat numele în *Web3D* și a început să lucreze la succesorul *VRML -X3D*.

VRML a provocat mult interes, dar nu a fost niciodată utilizat pe scară largă. Un motiv pentru acest lucru poate fi fost lipsa de lățime de bandă disponibilă. Experimentarea *VRML* a fost în primul rând în domeniul educației și cercetării.

VRML este util pentru o varietate de aplicații, inclusiv pentru vizualizarea datelor, analize financiare, Entertainment, educație, proiectare asistată de calculator, mall-uri virtuale, interfețe la informații, vizualizare științifică.

Scopul limbajului *VRML - Virtual Reality Modeling Language* a fost crearea unui standard prin care să se facă tranziția de la interfața bidimensională a Web-ului la una tridimensională și care să poată fi și manipulată dinamic.

VRML permite descrierea de obiecte 3D și combinarea lor în scene și lumi virtuale. El poate fi folosit la crearea de lumi interactive, care să conțină: imagini, animații, clipuri audio sau video. Într-o lume *VRML* utilizatorul și nu computerul are controlul. Obiectele dintr-o scenă pot interacționa unele cu altele prin evenimente sau pot interacționa cu evenimente utilizator. Lumile pot fi legate între ele sau cu alte documente

HTML. O lume creată poate fi distribuită pe Internet, văzută și explorată interactiv de mai mulți vizitatori în același timp.

Pentru a putea vizualiza o lume VRML, browser-ul de Internet (Internet Explorer, Netscape Navigator, HotJava, Opera etc.) are nevoie de un *VRML browser* (sub forma unui plugin de cele mai multe ori) ce permite utilizatorului (numit *vizitator* în VRML) să navigheze în lumea virtuală, să obțină diferite vederi ale scenei în funcție de punctul de vedere ales, să interacționeze cu elementele scenei etc.

4. Conversia documentelor din format tradițional în format electronic

Pe plan mondial, proiectele de digitizare atrag tot mai mulți cercetători, ingineri, specialiști în calculatoare, în tehnologia informației și a comunicării. Interesul mare în această direcție se explică prin fondurile substanțiale care au început să fie investite în acest tip de cercetare:

- Colorado Digitization Project a fost inițiat în 1998 și a reunit eforturile bibliotecilor, muzeelor, societăților istorice și arhivelor din Colorado pentru a crește accesul utilizatorilor la colecțiile speciale și resursele unice deținute de aceste instituții. Proiectul presupune un catalog comun de metadate și a dezvoltat instrumente speciale pentru creatorii de baze de metadate.

- Conferința IFLA din anul 2003, ținută în Turcia, a avut ca temă de dezbatere subiectul digitizării: „Introducing Digitization into Turkish Libraries: Current Attitudes and the Way Forward”.

- Un alt proiect important și complex de digitizare a fost lansat în Spania. AGID - Archivo General de Indias - a inițiat un proiect deosebit de ambițios pentru digitizarea colecțiilor sale. Mai mult de un milion de pagini de documente privind istoria Spaniei au fost digitizate.

- Uniunea Europeană a lansat, de asemenea, foarte multe proiecte de standardizare a cataloagelor de publicații, de digitizare a publicațiilor și de creare a portalurilor de acces multilingve la aceste biblioteci virtuale.

Noile tehnici de digitizare, recent dezvoltate, combină în mod optim robotica, electronica digitală și tehnica de calcul prin crearea liniei de digitizare care oferă în premieră o soluție complet integrată. Cele două avantaje majore ale liniei de digitizare sunt productivitatea foarte ridicată și menținerea integrității fizice a cărții supuse procesului de digitizare. Manipularea documentelor, transferul și digitizarea lor sunt automate, asigurându-se prezervarea fidelă a conținutului informației precum și a formei originale a documentului indiferent de tipul acestuia (cărți, reviste, colecții de ziare).

A. În octombrie 2002, firma elvețiană 4Digitalbooks ASSY, cea care a conceput DIGITIZING LINE, a instalat un prim echipament de acest tip la Green Library, una dintre bibliotecile Universității Stanford, California. Rezultat direct al unui prototip pentru formatul maxim A4, prima versiune industrială a lui DIGITIZING LINE permite scanarea complet automată a volumelor sau a ziarelor până la formatul A2. Echipamentul realizează în premieră mondială întoarcerea automată a paginii pentru întreaga gamă de formate de pagină.

Bibliotecile Universității Stanford au în total aproximativ opt milioane de cărți digitizate printre care toate titlurile publicate de Stanford University Press, inclusiv conservarea lucrărilor începând cu 1923. În prezent, sistemul DIGITIZING LINE instalat la Stanford funcționează la capacitate maximă pentru a scana câteva mii de cărți anual.

Proiectul prin care s-a dezvoltat linia de digitizare la Biblioteca Universității din Southampton se numește BOPCRIS. Una dintre principalele lui activități este scanarea documentelor parlamentare, rare și istorice, folosind DIGITIZING LINE, pentru prima dată în iunie 2004 în Marea Britanie și a doua oară în lume. Linia de digitizare scanează 600 de pagini / oră, ceea ce a făcut posibil accesul la documentele istorice, inaccesibile utilizatorilor până în acel moment. (DigiBook, nr. 9, December 2004).

În România, guvernul a început să acorde interes tot mai mare acțiunilor de digitizare, în special a materialelor vechi, pentru a oferi acces mai mare și pentru a le păstra intacte, în același timp. Câteva exemple prezentate în cele ce urmează arată stadiul actual al digitizării în țara noastră:

ANBPR (Asociația Națională a Bibliotecilor Publice) - Comisia pentru digitizarea documentelor: propune un proiect de digitizare ca o metodă de păstrare și conservare a documentelor în original;

CIMEC (Centrul Institutului de Memorie Culturală): lansează ideea digitizării de fotografii și crearea unei baze de date;

Biblioteca Universității „Transilvania” a inițiat programul „Leonardo da Vinci” în anul 2003. Proiectul „*Abordări inovatoare în cadrul bibliotecii universitare - formarea de bibliotecari și specialiști specializați în comunicare electronică, web-design, transmisie de informație digitală, capabili să conceapă produse și servicii utilizând ICT*” din anul 2003 s-a finalizat cu realizarea de produse și servicii de o calitate deosebită, rezultate diseminate în cadrul universităților și instituțiilor participante la proiect: S-au obținut documente digitizate diseminate folosind două instrumente importante: catalog Web și catalog pe CD.

Serviciile oferite se referă în primul rând la punerea la dispoziția publicului cititor a unei varietăți de resurse originale, incluzând cărți, documente, hărți, atlase, toate făcând parte din fondul rar și vechi și, în același timp foarte valoros. Devin accesibile informații complete referitoare la fondul rar: titlu, scanarea copertei și a paginii de titlu, descriere bibliografică, clasificare zecimală universală (CZU) și cuvinte cheie.

Un alt proiect „Leonardo da Vinci” a demarat în anul 2005, reunind instituții de renume din Brașov, Prima Școala Românească, Arhivele Naționale, filiala Brașov, și Biblioteca Județeană Brașov și având ca obiectiv crearea unui *Centru de Pregătire* prin intermediul căruia să se asigure pregătirea profesională continuă în cadrul organizațiilor, universităților, îmbunătățirea serviciilor publice și relațiilor interuniversitare, interorganizaționale, legăturile dintre universități și instituții sau alte companii. Un alt pas important a fost făcut prin caștigarea grantului CNCISIS în anul 2003, „*Cercetari privind metodele moderne de conservare a colecțiilor, cu aplicații în managementul bibliotecii digitale*”. În acest proiect, echipa de cercetare a testat și cercetat o nouă bază teoretică și experimentală pentru proiectare și realizare de sisteme tehnice de conservare și arhivare a documentelor tradiționale și a celor electronice din biblioteci.

În toate aceste proiecte românești soluția de digitizare aleasă a fost ***DIGITIZING LINE***, produsul realizat în Elveția și care a fost deja utilizat în Universități din întreaga lume. El realizează scanarea documentelor (cărți, reviste, ziare) cu ajutorul unei tehnologii foarte avansate, întoarcerea paginilor și stocarea informațiilor fără afectarea calității și integrității documentelor.

Posibilitatea de a folosi o gamă largă de dispozitive scanner este una din caracteristicile importante ale unui sistem de digitizare bun.

4.1. Metode de digitizare

4.1.1 Cerințe în procesul de digitizare

În acest capitol vom face o scurtă prezentare a diferitelor modalități în care se poate realiza achiziția imaginilor în formă digitală, plecând de la o carte. Unele variante au fost încercate în mod concret, rezultatele confirmând așteptările noastre.

O primă problemă care apare este cea a manipulării cărților în timpul procesului de digitizare, principalele probleme întâlnite fiind gradul ridicat de degradare, respectiv modul de legare al acestor cărți, de obicei destul de strâns, ceea ce împiedică deschiderea perfect dreaptă (fig. 4.1.).

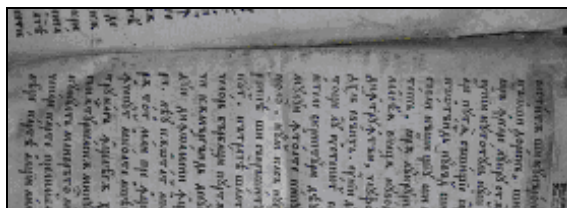


Fig. 4.1. Modul de legare al cărților

Pentru a evita deteriorarea cărților va trebui deci să luăm în considerare faptul că acestea nu vor putea fi manipulate în aceleași condiții ca documente nelegate.

În general există mai multe metode de a realiza digitizarea unei document:

1. Folosirea unui scanner ;
2. Folosirea unui aparat foto digital (camere planetare) ;
3. Scanarea filmelor folosite în fotografia clasică.

4.1.2 Scanare

Această variantă are ca avantaj incontestabil costul și rezoluțiile ridicate de scanare atinse, dar există o serie de limitări majore:

- Pentru a obține o imagine decentă, cartea trebuie să fie amplasată cu fața în jos pe sticlă, apoi ridicată, trebuie dată pagina și așezată din nou. Aceste operații repetate pot deteriora o carte.

- Cartea trebuie presată puternic la cotor pentru a face paginile cât mai paralele cu dispozitivul de scanare.

- Chiar și prin apăsare puternică, unele cuvinte sau imagini vor fi distorsionate, iar spre cotor va apărea o nuanță mai închisă.

- În cazul unor documente speciale, dimensiunea standard de scanare (unde va pe la A4) poate să nu ajungă.

4.1.3 Fotografiere digitală

O soluție acceptabilă din punct de vedere al costurilor, dar cu pierdere în rezoluție, comparativ cu soluția scanării. Avantajul major, poate chiar hotărâtor, este faptul că operațiile asupra cărții de digitizat sunt mai puțin agresive, fiind nevoie de simpla așezare și deschidere a cărții, pentru fiecare pagină fiind necesară doar operația de întoarcere a paginii. Cel mai important lucru este cel că nu trebuie să apăsăm cartea pentru a obține rezultate bune.

În scopul prelucrărilor ulterioare (de exemplu OCR), avem nevoie de anumite rezoluții minime. În următorul tabel avem câteva valori utilizate în descrierea aparatelor digitale și conversia acestora în alte unități de măsură. (tabel 5)

Tabel 5. Caracteristici tipice ale aparatelor digitale

Megapixeli	Dimensiuni (pixeli x pixeli)	Dimensiune pe disk (Mb)	Dimensiune in inch la 300 dpi
0.3	640 x 480	0.9	2.1 x 1.6
2.1	1792 x 1200	6.15	6 x 4
3.2	2160 x 1440	8.9	7.2 x 4.8
4.3	2400 x 1800	12.4	8 x 6
5.1	2608 x 1952	14.6	8.7 x 6.5
8	3264 x 2448	22.9	10.9 x 7.3

Un alt criteriu ce merită luat în calcul în cazul fotografierii digitale îl reprezintă numărul de puncte de focalizare pe care îl are camera digitală. Un număr mai mare de puncte de focalizare va asigura o calitate mai buna a imaginii, chiar dacă pagina nu este perfect întinsă, caz frecvent în cazul cărților cu număr mare de pagini unde manipularea poate fi dificilă.

Varianta cea mai bună ar fi folosirea unei camere planetare, aparat special construit pentru digitizarea cărților, dar cu costuri mult mai mari, având un preț orientativ de 15.000 dolari. Acest aparat combină avantajele camerei digitale și cele ale scannerului.

4.1.4 Scanarea fotografiilor analogice

O soluție extremă ar fi scanarea filmelor realizate prin fotografiere clasică. Pentru această avem nevoie de un film cu sensibilitate mare, un aparat analogic și un scanner care permite scanarea negativelor.

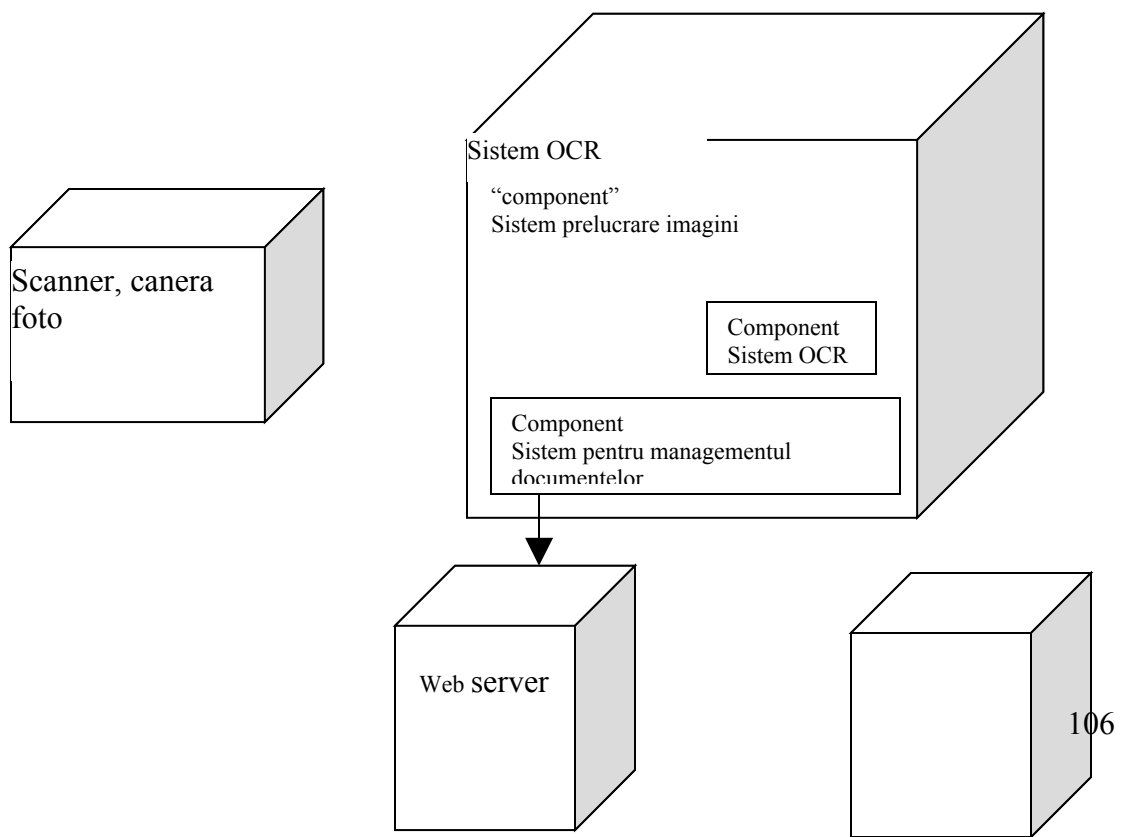
4.2 Arhitectura unui sistem de digitizare

Un asemenea sistem are ca scop îmbunătățirea procesului de recunoaștere de caractere prin utilizarea unui sistem hardware – software adecvat acestei activități și prin definirea unor procese clare de utilizare a acestuia.

Digitalizarea imaginilor se face cu ajutorul unei camere foto digitale sau a unui scanner special pentru a nu fi necesară deschiderea completă a cărților, astfel evitându-se deteriorarea acestora.

Sistemul este compus dintr-o componentă de achiziție de imagini (scanner sau cameră foto digitală), un sistem software de preprocesare a imaginilor, un sistem software de recunoaștere a caracterelor și un sistem de înmagazinare a informațiilor acumulate în diversele faze ale procesului.

Documentele rezultate pot fi apoi prezentate atât utilizatorilor web cât și utilizatorilor locali prin intermediul unor module de prezentare personalizate în funcție de cerințe. În figura următoare se pot vedea componentele importante ale sistemului și legăturile dintre acestea. Săgețile prezintă curgerea informațiilor în interiorul sistemului nu dependențe între componente (fig. 4.2.).



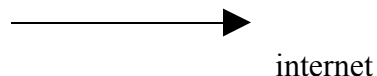


Fig.4.2. Schema bloc a sistemului de digitizare

Trebuie luat în considerare faptul că scopul final este o prezentare a documentelor unor utilizatori interesați astfel că ieșirile din acest sistem ar trebui să fie documente care să permită căutarea în textul acestora pe cât posibil fără a deteriora informația vizuala din documentele originale.

Documentele create trebuie să aibă următoarele caracteristici:

- trebuie să conțină atât imagini cât și text;
- să existe posibilitatea de căutare;
- să poată fi citite de către clienți web uzuali;
- să aibă conceptul de pagină;
- să fie de dimensiuni reduse ;
- să fie suportate de către sistemele OCR existente (ca documente de ieșire).

4.3 Moduri de livrare a conținutului digital

Pentru a livra conținutul (imaginile pe care s-a făcut OCR) propunem ca soluție documentele Acrobat. Acest tip de documente conțin mai multe nivele informaționale. Nivelul text se poate pune atât deasupra nivelului imagine cât și sub acesta. Cele mai utile modalități de prezentare a conținutului în acest caz sunt: text sub imagine și respectiv text peste imagine.

4.3.1 Text sub imagine

În literatură acest mod de prezentare se numește “text under the image”. Adobe a creat această funcționalitate pentru a pune la dispoziție o funcționalitate de căutare în documentele cu conținut “sensibil” în care recunoașterea nu este perfectă și în care nu ne permitem ca utilizatorul final să vadă greșelile dar totuși existând posibilități de selecție / căutare a textului de sub imagini, selecția reflectându-se asupra imaginii de deasupra.

4.3.2 Text peste imagine

În literatură este cunoscut sub numele de “text over the image”. Textul original de pe imaginile pe care s-a făcut OCR se șterge automat și este înlocuit de textul recunoscut.

4.4. Alegerea unui format

Prin analiza criterială am stabilit importanta caracteristicilor și am luat în considerare mai multe tipuri de documente existente inclusiv un tip de document proprietar care ar putea fi dezvoltat special pentru acest proiect (tabel 6).

În urma analizei s-a ajuns la concluzia că dezvoltarea unui format proprietar nu este cea mai bună soluție, documentul PDF câștigând:

Tabel 6 : Caracteristici ale diferitelor formate de documente

Caracteristica	Importanță caracteristică	. TIFF	DOC	PDF	Custom format
Image + text	3	1	1	1	1
Căutare	4	0	0.5	1	1
Compatibilitate cu clienții existenți	5	0.75	1	1	0
Paginabile	2	1	0.75	1	1
Dimensiune	1	0.75	0.6	0.8	1
Suportabilitate	6	1	1	1	0
Punctaj acumulat		12.5	18.1	20.8	10

În plus față de celelalte documente, formatul PDF permite așezarea textului recunoscut sub imaginea inițială (text under the image) astfel că utilizatorul va vedea imaginile dar va căuta în text.

Acest lucru asigură o calitate ridicată a vizualizării chiar și în cazul în care recunoașterea a fost de o calitate mai scăzută. Documentele PDF permit și înmagazinarea de metainformații în acestea.

Mecanismele de protecție ale documentelor Pdf nu sunt de neglijat fiind posibil să protejăm un document la salvare, copiere, tipărire, copy-paste, etc .

Sistemul trebuie să fie gândit astfel ca toate documentele intermediare să fie păstrate, în cazul unor erori nemaifiind necesară reluarea întregului proces (de la scanare la salvare PDF) ci doar părți ale acestuia.

4.5. Tehnici de conversie prin Recunoașterea Optică a Caracterelor (Optical Character Recognition-OCR)

4.5.1. Scurt istoric al conversiei documentelor din format tradițional prin scanare și Recunoașterea Optică a Caracterelor (OCR)

Recunoașterea textului din imagini a fost un subiect mult discutat de-a lungul timpului. „Recunoașterea optică a caracterelor (OCR) transformă imagini de text, cum ar fi documentele scanate, în caractere de text. Cunoscută și sub numele de recunoașterea textului, OCR face posibilă editarea și reutilizarea textului conținut de imaginile scanate. OCR utilizează o formă de inteligență artificială, cunoscută sub numele de recunoașterea modelului, pentru identificarea individuală a caracterelor unui text dintr-o pagină, inclusiv semnele de punctuație, spațiile și sfârșitul de linie”. (<http://office.microsoft.com/ro-ro/help/HP030812551048.aspx>). De la începuturile din 1950 tehnica a fost în permanentă îmbunătățită. Script-urile recognoscibile au fost la început numerele arabe și apoi alfabetele latine, japoneze, chineze. Multe tipuri diferite de formate pe hârtie pot fi citite astăzi prin OCR-izare. Tehnicile de recunoaștere a caracterelor au cunoscut perioade diferite de dezvoltare. Se identifică două momente în care acestea au avut de suferit ca utilitate și dezvoltare. Primul moment a fost în anii '80, atunci când au apărut programele de birotică. Acestea facilitau crearea documentelor direct în format digital (documente "born digital"). Al doilea moment a fost implementarea crescândă a noilor tehnologii și utilizarea Internetului. Se reevaluează importanța manuscriselor și trecerea lor în format digital. Dispozitive mobile mobile cu microcamere au acum incorporate unități de procesare capabile de recunoaștere în timp real al caracterelor.

OCR a apărut în 1950 în Statele Unite ale Americii, în aceeași perioadă în care apărea calculatorul UNIVAC. În anii 1960 IBM crează propriul program de recunoaștere capabil să citească numere tipărite și scrise.. Tot în acest an, s-au mecanizat operațiunile poștale, astfel scrisorile erau sortate cu ajutorul unor dispozitive mecanice cu OCR. În 1974, Ray Kurzweil a dezvoltat un program capabil de a recunoaște caracterele tipărite în orice font. În anii 1980 dispozitivele de recunoaștere și-au redus dimensiunile datorită progreselor din sfera semiconductorilor și a microprocesoarelor.

4.5.2. Starea actuală a tehnologiei OCR

Recunoașterea precisă a fontului latin, text scris la mașină nu este considerat o problema rezolvată în aplicațiile unde imaginile clare sunt puse la dispoziție prin scanarea documentelor printate.

Rata preciziei depășește 99%; acuratețea totală poate fi atinsă doar prin reverificare umană. Alte arii cum ar fi cele care includ recunoașterea scrisului de mână și a celui printat în alte fonturi (scripts), în special cele cu un număr mare de caractere, sunt în continuare subiectul cercetărilor în domeniu.

Rata preciziei poate fi măsurată în mai multe feluri, precum și modul în care acestea sunt măsurate pot afecta foarte mult rata raportată de precizie. De exemplu, dacă nu este folosit un dicționar pentru a corecta cuvintele nonexistente găsite de către soft, o marja de eroare de 1% la litere (acuratețe de 99%) poate duce la o marja de eroare de 5% (acuratețe de 95%) sau mai mult în cazul în care fiecare cuvânt cu o litera greșită este luat ca și greșit.

Recunoașterea on-line a caracterelor este deseori confundată cu recunoașterea optică a caracterelor. OCR este un sistem de recunoaștere de caractere **off-line**, unde sistemul recunoaște formele statice a caracterelor, în timp ce recunoașterea **on-line** a caracterelor implică recunoașterea mișcării dinamice a scrisului de mână. De exemplu, recunoașterea on-line, cum ar fi recunoașterea gesturilor în SO Penpoint sau în Tablet PC, poate preciza dacă o linie orizontală a fost desenată dinspre dreapta spre stânga sau invers. Recunoașterea on-line a caracterelor este, de asemenea, echivalentă și cu alți termeni, cum ar fi: recunoașterea dinamică a caracterelor, recunoașterea în timp real a caracterelor și recunoașterea inteligentă a caracterelor denumită și ICR.

Sistemele dinamice de recunoaștere on-line au devenit cunoscute ca produse comerciale în ultimii ani. Printre acestea se numără și dispozitivele periferice pentru asistență personală digitală asemenea celor care folosesc SO Palm. Corporația Apple a patentat acest produs. Algoritmii folosiți în aceste dispozitive au avantajul faptului că ordinea, viteza și direcția liniilor de segment individuale la introducere sunt cunoscute.

De asemenea, utilizatorul este obligat să folosească doar forme specifice. Aceste metode nu pot fi folosite în softuri care scanează documente de hârtie, astfel recunoașterea cu acuratețe a documentelor scrise de mână este încă o problema mare.

Rata acurateții de 80% până la 90% a fișierelor cu scrisul de mână lizibil poate fi atinsă, dar cu o asemenea acuratețe tot mai apar zeci de greșeli pe pagină, astfel aceasta tehnologie este folositoare doar în aplicații limitate.

Recunoașterea textului cursiv este o zona activă de cercetare, cu rate de recunoaștere chiar mai mici decât cea a recunoașterii scrisului de mână. Rate mai mari de recunoaștere a textului nu vom putea obține fără ajutorul informațiilor gramaticale. De exemplu, recunoașterea cuvintelor întregi folosind un dicționar este mai ușoară decât încercarea de analiză individuală a caracterelor.

Cunoașterea gramaticii limbii textului scanat poate de asemenea să ajute, de exemplu, să determine dacă un cuvânt este verb sau substantiv pentru o mai mare acuratețe. Formele individuale de caractere cursive pur și simplu nu conțin informații suficiente pentru a recunoaște cu acuratețe (mai mult de 98%) toate caracterele scrisului de mână. Este necesar a înțelege că tehnologia OCR este o tehnologie de bază de asemenea folosită de aplicațiile de scanarea avansată. Din aceste motive, o soluție de scanare avansată poate fi unică, patentată și protejată de drepturile de autor, deși este bazată pe tehnologia OCR de bază.

Recunoașterea optică a caracterelor se realizează în doi pași:

-utilizarea unui dispozitiv pentru scanarea informației tipărite într-un calculator ca imagine bit-map;

-aplicarea algoritmilor de recunoaștere a caracterelor pentru crearea fișierului text.

Acuratețea atinsă în recunoașterea caracterelor din imagini este cel mai important factor care determină eficacitatea și fezabilitatea unui produs OCR.

Software-ul OCR este proiectat pentru a asigura conversia documentelor scrise olograf sau tipărite, digitizate prin scanare, într-o formă care se pretează procesării computaționale. Fișierele text în reprezentare ACSII sau Unicode sunt produse ale perogramelor de OCR-izare. Sistemele OCR își au originile în recunoașterea "pattern-urilor" și inteligența artificială. Recunoașterea vizuală a caracterelor utilizând tehnici ca oglinzi și lentile și recunoașterea digitală a caracterelor utilizând scanere și algoritmi specifici se regăsesc în tehnicile OCR.

4.5.3. Tehnologii curente, produse software pentru OCR (tabel 7)

Tabel 7. Tehnologii curente, produse software pentru OCR

Nume	Licența	Sistem de operare	Descriere
ExperVision Typereader OpenRTX	Comercial	Windows, Mac OS X, Unix, Linux, OS/2	ExperVision Inc. a fost fondat în 1987, tehnologia și produsul lor OCR. Au luat cele mai mari note în testarea independentă făcută de UNLV în anii în care a participat
ABBYY FineReader OCR		Windows, Mac OS X	Pentru a lucra cu interfețe locale este necesar suportul lingvistic corespunzător.
AnyDoc Software OCR for AnyDoc	Comercial	Windows	Lucrează cu structuri, semistrukturi și documente nestructurate
OmniPage	Comercial	Windows, Mac OS	Produs de Nuance Communications
Readiris Windows,	Comercial	Mac OS X	Produs de I.R.I.S. Group of Belgium. Ediții Asian și Middle Eastern
CuneiForm	BSD variant	Windows, Linux, BSD, MacOSX.	Enterprise-class system, multi language, poate salva textul formatat și tabele de recunoaștere complexe a oricărei structuri
Puma.NET	BSD	Windows	.NET OCR SDK bazat pe tehnologie
CVISION Technologies, Inc. PdfCompressor and Maestro Recognition Server	Comercial	Windows	Rapidă, precisă, volum mare de date.
GOCR	GPL	diverse (open source)	Dezvoltare de început.
Microsoft Office Document Imaging	Comercial	Windows	Folosește motorul ScanSoft OCR .
Microsoft Office OneNote 2007 1	Comercia	Windows	
NEOPTec DATAScan	Comercial	Windows	Soft pentru procesare automată a cererilor și

.			chestionarelor.
NovoDynamics VERUS	Comercial		Produs specializat în limbile din orientul mijlociu
Ocrad	GPL	Unix-like, OS/2	
Brainware	Comercial	Windows	Extracții și procesări de date din documente în orice sistem backend; printre documentele cunoscute se număra: chitanțe, declarații
HOOCR	GPL	Linux	Hebrew OCR
OCROPUS	Apache	Linux	Poate utiliza Tesseract
PDF OCR X	Comercial	Mac OS X	Utilitate drag and drop care poate converti fișiere PDF în fișiere text folosind OCR. Folosește Tesseract
ReadSoft	Comercial	Windows	Scanează și clasifică documente oficiale cum ar fi: chitanțe, facturi
Alt-N Technologies' RelayFax Network Fax Manager	Comercial	Windows	Utilitar OCR multi-lingvistic care convertește documente fax în documente editabile (.doc, .pdf,...) în limbi diferite.
Scantron Cognition	Comercial	Windows	Pentru lucrul cu interfețe locale, suport lingvistic
SimpleOCR	Freeware și Comercial	Windows	
SmartScore	Comercial	Windows, Mac OS	Pentru note muzicale.
Tesseract	Apache	Windows, Mac OS X, Linux	Creat de către Hewlett-Packard; sub dezvoltare curentă de către Google

4.5.4. Suportul lingvistic al produselor pentru OCR

Tabel:7 Suportul lingvistic al produselor OCR

Nume Ultima versiune Anul lansării	Limbi recunoscute
ExperVisionTypeReader & OpenRTK	English, French, German, Italian, Spanish, Portuguese, Danish, Dutch, Swedish,

7.0 2007	Norwegian, Hungarian, Polish, Simplified Chinese, Traditional Chinese, Russian, Finnish și Polynesian
ABBYYFineReader OCR 10.0 2009	Abkhaz, Adyghian, Afrikaans, Agul, Albanian, Altai, Armenian (Eastern, Western, Grabar), Avar, Aymara, Azerbaijani (Cyrillic), Azerbaijani (Latin), Bashkir, Basic, Basque, Belarusian, Bemba, Blackfoot, Breton, Bugotu, Bulgarian, Buryat, C/C++, COBOL, Catalan, Cebuano, Chamorro, Chechen, Chinese Simplified, Chinese Traditional, Chukchee, Chuvash, Corsican, Crimean Tatar, Croatian, Crow, Czech, Dakota, Danish, Dargwa, Dungan, Dutch (Netherlands and Belgium), English, Eskimo (Cyrillic), Eskimo (Latin), Esperanto, Estonian, Even, Evenki, Faroese, Fijian, Finnish, Fortran, French, Frisian, Friulian, Gagauz, Galician, Ganda, German (Luxemburg), German, Greek, Guarani, Hani, Hausa, Hawaiian, Hebrew, Hungarian, Icelandic, Ido, Indonesian, Ingush, Interlingua, Irish, Italian, JAVA, Japanese, Jingpo, Kabardian, Kalmyk, Karachay-balkar, Karakalpak, Kasub, Kawa, Kazakh, Khakass, Khanty, Kikuyu, Kirghiz, Kongo, Koryak, Kpelle, Kumyk, Kurdish, Lak, Latin, Latvian, Lezgi, Lithuanian, Luba, Macedonian, Malagasy, Malay, Malinke, Maltese, Mansy, Maori, Mari, Maya, Miao, Minangkabau, Mohawk, Moldavian, Mongol, Mordvin, Nahuatl, Nenets, Nivkh, Nogay, Norwegian (nynorsk and bokmål), Nyanja, Occidental, Ojibway, Ossetian, Papiamento, Pascal, Polish, Portuguese (Portugal and Brazil), Provençal, Quechua, Rhaeto-romanic, Romanian, Romany, Rundi, Russian, Russian (old spelling), Rwanda, Sami (Lappish), Samoan, Scottish Gaelic, Selkup, Serbian (Cyrillic), Serbian (Latin), Shona, Simple chemical formulas, Slovak, Slovenian, Somali, Sorbian, Sotho, Spanish, Sunda, Swahili, Swazi, Swedish, Tabasaran, Tagalog, Tahitian, Tajik, Tatar, Thai, Tok Pisin, Tongan, Tswana, Tun, Turkish, Turkmen, Tuvian, Udmurt, Uighur (Cyrillic), Uighur (Latin), Ukrainian, Uzbek (Cyrillic), Uzbek (Latin), Welsh, Wolof, Xhosa,

	Yakut, Zapotec, Zulu
OmniPage 17 2009	Afrikaans, Albanian, Aymara, Basque, Bemba, Blackfoot, Breton, Bugotu, Bulgarian, Byelorussian, Catalan, Chamorro, Chechen, Corsican, Croatian, Crow, Czech, Danish, Dutch, English, Esperanto, Estonian, Faroese, Fijian, Finnish, French, Frisian, Friulian, Gaelic (Irish), Gaelic (Scottish), Galician, Ganda/Luganda, German, Greek, Guarani, Hani, Hawaiian, Hungarian, Icelandic, Ido, Indonesian, Interlingua, Italian, Inuit, Kabardian, Kasub, Kawa, Kikuyu, Kongo, Kpelle, Kurdish, Latin, Latvian, Lithuanian, Luba, Luxembourgian, Macedonian, Malagasy, Malay, Malinke, Maltese, Maori, Mayan, Miao, Minankabaw, Mohawk, Moldavian, Nahuatl, Norwegian, Nyanja, Occidental, Ojibway, Papiamento, Pidgin English, Polish, Portuguese (Brazilian), Portuguese, Provencal, Quechua, Rhaetic, Romanian, Romany, Ruanda, Rundi, Russian, Sami Lule, Sami Northern, Sami Southern, Sami, Samoan, Sardinian, Serbian (Cyrillic), Serbian (Latin), Shona, Sioux, Slovak, Slovenian, Somali, Sorbian, Sotho, Spanish, Sundanese, Swahili, Swazi, Swedish, Tagalog, Tahitian, Tinpo, Tongan,
Readiris 12 Pro & Corporate 2009	American English, British English, Afrikaans, Albanian, Aymara, Balinese, Basque, Bemba, Bikol, Bislama, Brazilian, Breton, Bulgarian, Byelorussian, Catalan, Cebuano, Chamorro, Corsican, Croatian, Czech, Danish, Dutch, Esperanto, Estonian, Faroese, Fijian, Finnish, French, Frisian, Friulian, Galician, Ganda, German, Greek, Greenlandic, Haitian (Creole), Hani, Hiligaynon, Hungarian, Icelandic, Ido, Ilocano, Indonesian, Interlingua, Irish (Gaelic), Italian, Javanese, Kapampangan, Kicongo, Kinyarwanda, Kurdish, Latin, Latvian, Lithuanian, Luxemburgh, Macedonian, Madurese, Malagasy, Malay, Maltese, Manx (Gaelic), Maori, Mayan, Minangkabau, Nahuatl, Norwegian, Numeric, Nyanja, Nynorsk, Occitan, Pidgin English, Polish, Portuguese, Quechua, Rhaeto-Roman, Romanian, Rundi, Russian, Samoan, Sardinian, Scottish (Gaelic), Serbian, Serbian (Latin), Shona, Slovak, Slovenian, Somali, Sotho, Spanish,

	Sundanese, Swahili, Swedish, Tagalog, Tahitian, Tok Pisin, Tonga, Tswana, Turkish, Ukrainian, Waray, Wolof, Xhosa, Zapotec, Zulu, Bulgarian - English, Byelorussian - English, Greek - English, Macedonian - English, Russian - English, Serbian - English, Ukrainian - English + Moldovan, Bosnian (Cyrillic
Readiris 12 Pro & Corporate Middle-East 2009	Arabic, Farsi și Hebrew
Readiris 12 Pro & Corporate Asian 2009	Simplified Chinese, Traditional Chinese, Japanese și Korean
CuneiForm 12 2007	English, German, Croatian, Polish, Danish, Portuguese, Dutch, Digits, Czech, French, Romanian, Hungarian, Bulgarian, Slovenian, Lettish, Lithuanian, Estonian, Turkish, Russian, Swedish, Spanish, Italian, Russian-English (mixed), Ukrainian, Serbian
Microsoft Office Document Imaging Office 2007 2007	Accesul la diferite limbi este legat de instalarea unor componente MS Office
NEOPTec DATA-SCAN 5.7 2009	French, Spanish, English.
NovoDynamics VERUS Middle East Professional 2005	Arabic, Persian (Farsi, Dari), Pashto, Urdu, inclusiv English and French
NovoDynamics VERUS Asia Professional 2009	Chineza simplificată și tradițională , limbile Korean și Russian, incluzând English
HOCR 0.10.13	Hebrew
OCROPUS 0.3.1, 08	Toate limbile suportate de Tesseract prin plug-in-uri, și suportă Latin script și English în mod nativ
ReadSoft Europene,	Caractere Chineza simplificată și tradițională, caractere Korean și Japanese

SimpleOCR 3.5 2008	Engleză, franceză
Tesseract 2.03 2008	Poate recunoaște 6 limbi, compatibil UTF8, are suport de antrenare

4.5.5. Particularități ale sistemelor OCR actuale

- *ExperVision*

Sisteme de operare: Windows, Mac OS X, Unix, Linux, OS/2

A câștigat cele mai mari note la testele de performanță făcute de UNLV tehnologiilor OCR.

OpenRTK este o tehnologie pentru clienții OEM sau care au nevoie să integreze OCR-ul în aplicațiile lor (exemplu DIM, FPS).

TypeReader 2008 este un program de recunoaștere a caracterelor creat în special pentru segmentul corporate. Este adresat celor ce au nevoie de transformarea a sute de pagini în documente electronice, dar numai dacă ai nevoie de a face acest lucru rapid și nu te interesează prea mult calitatea. Viteza este de aproximativ 3 ori mai bună decât cea a lui Abbyy Fine Reader Professional Edition și de 4 ori mai bună decât a lui OmniPage Professional. La fel ca și OmniPage Professional, TypeReader include caracteristica "watched folder". Aceasta presupune posibilitatea de a aplica automat procesul de OCR-izare tuturor documentelor ce sunt plasat într-un anumit director din rețea. Viteza de recunoaștere a caracterelor este absolut uimitoare, iar prezența unei astfel de caracteristici face și mai mare viteza de recunoaștere a caracterelor care era mare. Totuși, unele asocieri de litere nu sunt recunoscute corespunzător (exemplu "rt" sau "tr").

- *ABBYY*

Avantajele aplicației Fine Reader OCR 8.0 constau în procesarea facilă a documentelor în format pdf, recunoașterea pozelor realizate cu camera digitală și recunoașterea optică a caracterelor printr-o singură apăsare de buton. Asigurând o capacitate semnificativ îmbunătățită de convertire a imaginilor documentelor în format HTML, FineReader 10 simplifică citirea și publicarea cărților electronice. Are funcții de

analiza și păstrare a formatului paginii. Oferă o optimizare a dimensiunii documentelor PDF rezultate, reducându-le de până la 10 ori, păstrând însă calitatea vizuală.

Recognition (ADRTTM) analizează documentul ca o entitate singulară, păstrând toate elementele din structura documentului original, inclusiv textul, tabelele, antetul / subsolul și numerotarea paginilor. Are dicționare în peste 30 de limbi.

- ***AnyDoc***

Dezvoltă tehnologii pentru a procesa documente structurate, semi-structurate și nestructurate, de clasificare. Documente structurate, cum ar fi o cerere de credit, folosesc un anumit șablon.

- ***OmniPage***

Convertește imagini, cum ar fi documentele scanate și fișierele PDF, în formate de fișiere folosite de aplicații cum ar fi Microsoft Word, Excel, Adobe Acrobat, sau de fișiere HTML.

- ***Readiris***

În 2006, a fost selectat pentru a fi utilizat în Adobe Acrobat 8.0. Softul se află în concurență cu Expervision (TypeReader), ABBYY FineReader, OmniPage, GOCR și Tesseract.

- ***CuneiForm***

Cuneiform poate transforma texte în format RTF, HTML, text simplu, Microsoft Word sau Microsoft Excel. Poate recunoaște tabele complicate de orice structură. În 2008 Cognitive Technologies a lansat on-line servicii gratuite de recunoaștere pe penOCR.org.

- ***Puma.NET***

Poate încărca o singură pagină, o recunoaște și salvează rezultatul într-un fișier. Poate să returneze rezultatul recunoașterii într-un șir de caractere. Puma.NET nu poate să încarce pagini multiple și să salveze rezultatul într-un singur fișier. Nu oferă utilizatorului posibilitatea de a vedea în paralel originalul și varianta recunoscută.

- ***CVISION Technologies***

Oferă o calitate superioară pentru documentele colorate și procesează un volum mare de documente.

Localizează rapid un cuvânt într-un document de dimensiuni mari. În plus este optimizat pentru WEB și permite procesarea PDF-urilor.

- ***OCRopus***

OCRopus este un sistem de ultima oră pentru analiza documentelor și OCR, ce conține un analizator de layout-uri, un dispozitiv de recunoaștere a caracterelor, modelare statistică a limbajului natural și capacități multilingvistice.

Motorul OCRopus se bazează pe două proiecte de cercetare: un sistem de recunoaștere dezvoltat în anii 90' de biroul US Census și pe metode moderne, de mare performanță, de analiza a layouturilor.

Dezvoltarea OCRopus este sponsorizată de către Google. Inițial proiectul a fost gândit ca un ajutor pentru eforturile de conversie masivă a documentelor. Se dorește de asemenea a fi un excellent sistem OCR pentru multe alte aplicații.

- ***PDF OCR X***

PDF OCR este un program simplu care convertește un document PDF în document text. Acesta utilizează tehnologie OCR avansată prin care poate să extragă textul din PDF chiar dacă textul este conținut într-o imagine.

Avantaje:

1. Lucrează cu orice fel de PDF chiar dacă este un PDF scanat sau un PDF generat dintr-un anumit document;
2. Interfața drag and drop ușor de folosit;
3. Poate fi utilizat mai multe limbi: Engleză, franceză, germană și italiană;

- ***Alt-N Technologies- RelayFax Network Fax Manager***

RelayFax network fax server management software automatizează trimiterea, primirea și managementul faxurilor rețelei din desktopul tau, în timp ce integrează total funcționalitățile faxului cu sistemul de email existent - totul în timp ce furnizează trimiteri și primiri nelimitate de faxuri, fără nici o taxă lunară sau per-transmitere.

RelayFax se conectează la un număr nelimitat de POP mailboxes, la intervale programate și colectează mesajele de fax aflate în așteptare pe care apoi software-ul le trimite prin e-mail, fax sau printează după cum este definit în regulile configurabile ale fax-ului.

Avantaje:

- detectare automată a modemului;
- utilizarea mediului Terminal Services;
- programarea fax-urilor;
- **conversia fax-urilor în PDF sau PNG;**
- tehnologie modernă de trimitere a fax-urilor;

- ***Scantron***

Este o soluție de business care elimină în mod eficient procesele manuale care nu fac diferența între informațiile primite prin fax, mail sau internet.

Avantaje:

1. Reduce timpul și costul cheltuit pe introducerea manuală a datelor cu până la 90%;
2. Trimiterea reală și precisă a datelor și a documentelor;
3. Selecție variată de programe folosite (Visual Basic, C++, VB Script, .NET).

- ***SmartScore***

Acest software ajută să aveți muzică printată într-un mod foarte rapid în computer. Asigură recunoașterea oricărui format PDF.

- ***Tesseract***

Avantaje:

1. Poate identifica dacă textul este proporțional sau monospațiu;
2. Poate fi utilizat în mai multe limbi;
3. Dezavantaje:
4. Este utilizat doar pentru imaginile de format TIFF (asta doar dacă libtiff este instalat);
5. Poate fi utilizat pe Windows, Linux, Mac OS X, dar din cauza resurselor limitate

dezvoltatorii l-au testat doar pentru Windows si Linux.

- ***ReadSoft***

ReadSoft Documents este o aplicație complexă destinată procesării documentelor de tip structurat și semistructurat. Este o aplicație completă, adaptabilă oricăror cerințe de business sau volum de documente. Este capabilă de a captura toate formularele de procesat, indiferent de format și / sau design. Se poate integra ușor cu aplicațiile proprii (SAP, Oracle).

- ***Microsoft Office Document Imaging***

Se recomandă a se utiliza Microsoft Office Document Imaging atunci când se dorește:

- Scanarea documentelor de o singură pagină sau de mai multe pagini. De exemplu, se pot scana documente pe hârtie pentru arhivare și reciclarea copiilor pe hârtie;

- Citirea unui document scanat sau a unui fax, rapid și ușor, pe ecranul computerului. De exemplu, se poate citi în mod conectat un fax pe mai multe pagini;

- Completarea în mod conectat a unui formular simplu care a fost scanat sau primit ca fax. De exemplu, scanați un formular pe hârtie sau deschideți un formular primit ca fax, îl completați cu informațiile solicitate, inclusiv casetele de selectare, apoi returnați formularul prin poșta electronică ;

- Efectuarea recunoașterii optice a caracterelor (OCR) într-un document scanat sau într-un fax ;

De exemplu, după recunoașterea textului din documentul scanat sau din fax, se poate căuta un text specific sau se poate copia text în alt program.

- Copierea de text și imagini dintr-un document scanat sau dintr-un fax și lipirea lor în orice program Office. Exportul de text și imagini dintr-un document scanat sau dintr-un fax în Microsoft Word. De exemplu, dace se dorește copierea de date importante dintr-un fax sau dintr-un document scanat într-o foaie de lucru Microsoft Excel ;

- Căutarea de text într-un document scanat sau într-un fax ;

- Reorganizarea ordinii paginilor într-un document scanat sau fax cu mai multe pagini, la fel de ușor ca rearanjarea hârtiilor într-un dosar. De exemplu, eliminați pagina

de însoțire a unui fax sau îi adăugați pagini suplimentare înainte de a-l trimite la altă persoană ;

-Trimiterea de documente scanate către alții, prin poștă electronică sau ca fax, prin Internet ;

- Adnotarea unui document scanat sau a unui fax și partajarea lor cu altă persoană.

De exemplu, adăugarea de comentarii unui fax și returnarea lui la expeditor.

Microsoft Office Document Imaging are două componente — una de scanare și una pentru lucrul cu imagini.

Componenta de scanare permite scanarea documentelor pentru a le face disponibile în computer, utilizând orice scener instalat. Această componentă furnizează setări prestabilite de scanare pentru controlul scannerului utilizând setări optimizate pentru scopuri specifice. De exemplu, setarea prestabilită de scanare Alb-negru dă cele mai bune rezultate la OCR când scanați pagini de text pentru OCR, în timp ce setarea prestabilită de scanare Color este cea mai bună pentru a scana imagini colorate sau lucrări de artă. OCR este efectuat automat pe documente text imediat după scanare și se pot scana ușor pagini multiple într-un singur fișier.

Componenta de imagine facilitează vizualizarea pe ecran a documentelor scanate, rearanjează documentele multipagină, selectează și manevrează textul recunoscut, adnotează documentele scanate și fax-urile electronice și trimite documente către alții prin poștă electronică sau fax.

- **Fine Reader 7.0**

Deoarece există o multitudine de documente în format clasic, pe hârtie, care trebuie și ele aduse în digitală este foarte procesul de OCR

formă important

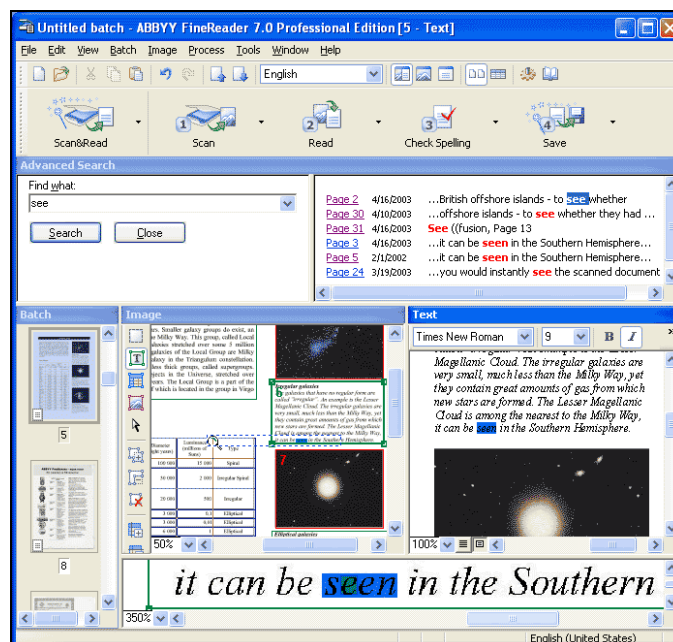


Fig.4.3. Utilitarul FineReader 7.0 de la ABBYY

(Optical Character Recognition), în cele ce urmează prezentând unul din cele mai de succes utilitare în acest sens - FineReader.

Datorită faptului că utilizează tehnologia IPA, FineReader permite recunoașterea foarte bună a conținutului. Renumit deja datorită replicărilor perfecte pentru proiecte complexe și pentru tabele și documente de calitate slabă, acum prezintă îmbunătățiri suplimentare. Sunt proiectate proceduri noi de filtrare 'Intelligent background' pentru filtrarea fundalului color, zgomotelor și altor interferențe și permite recunoașterea și mai corectă a unor documente care până acum constituiau o adevărată provocare pentru utilitarele de OCR-izare, precum: articole cu un fundal foarte colorat, texte tipărite peste imagini, sau urme lăsate de agrafe. Acum nu mai are importanță ce fel de documente se scanează, fie ele poze multi-coloană, texte încadrate, sau prețuri dintr-un ziar, pentru că FineReader asigură acuratețea necesară.

Se pot recunoaște texte scrise în 177 de limbi, inclusiv engleza, germana, franceza, spaniola, italiana, portugheza, olandeza, suedeza, finlandeza, rusa, ucraineana, bulgara, ceha, maghiara, poloneza, slovacă, malaeziana, indoneziana și altele. Verificarea ortografiei este disponibilă pentru 34 de limbi. Acest fapt simplifică procesul de introducere a datelor pentru oameni aparținând diferitelor țări și culturi.

Suportă exportul către aplicațiile uzuale de birou, inclusiv Microsoft Word, Microsoft Excel, Microsoft PowerPoint, Lotus Word Pro, Corel WordPerfect și Sun StarWriter. Textul recunoscut poate fi salvat într-o varietate de formate ale fișierelor, inclusiv PDF (cu șase opțiuni de salvare), HTML, Microsoft Word XML, DOC, RTF, XLS, PPT, DBF, CSV și TXT. Se pot converti cu ușurință documentele în format XML, utilizând Microsoft Word XML, un standard în devenire pentru utilizarea fișierelor. Acesta permite utilizatorilor să beneficieze de avantajele fișierelor în format XML,

putând fi manipulat și căutat prin utilizarea oricărui program care prelucrează standardul XML. Cu ajutorul suportului XML, FineReader 7.0 simplifică procesul de citire și editare a documentelor word. În momentul exportului documentelor în Microsoft Office Word 2003, FineReader deschide în mod automat imaginea inserată a documentului original, permițând utilizatorilor să editeze și să vadă în mod simultan documentele originale, eliminându-se astfel necesitatea de comutare între cele două aplicații.

Se pot deschide, citi și edita fișiere în PDF, care pot fi salvate în orice format sau poate exporta datele către aplicația favorită. Toate fișierele PDF create în FineReader sunt optimizate pentru publicarea pe Web. FineReader suportă exportul în oricare dintre cele patru formate standard ale fișierelor PDF (numai text și imagini, text suprapus imaginii, text sub imagine, și numai imagine), cu controale suplimentare. Aceasta oferă utilizatorilor șase opțiuni de salvare în PDF.

Este asigurat suportul pentru periferice multi-funcționale care combină funcționalitățile scannerului, ale imprimantei, fotocopiatorului și faxului, atât conectate la o stație de lucru, cât și în rețea. Setările speciale ale programului permit utilizatorilor să deschidă în mod automat imagini scanate, din orice locație a rețelei sau de la un server FTP și să le recunoască.

5. Elaborarea unui model conceptual pentru un sistem de informare și documentare

5.1. Depozite digitale

Resursele digitale aflate în permanență dezvoltare formează baza capitalului intelectual în cercetare. Acestea trebuie să persiste și să poată fi căutate, accesate și înțelese. Reutilizarea datelor (de către utilizatori din domenii diferite), trebuie să fie posibilă imediat ce acestea au fost create cât și pentru o perioadă mai mare de timp. Aceleași tehnici de preservare a datelor vor susține formarea depozitelor actuale cât și utilizarea lor ulterioară. Există riscul însă ca multe dintre datele științifice și tehnice să se piardă pentru generațiile viitoare dacă nu sunt prezervate corespunzător.

Cererile în domeniu sunt concentrate pe preservarea informațiilor digitale în știință, pornind de la datele primare până la analizarea publicațiilor finale rezultată din cercetare. Este important a dezvolta o hartă (roadmap) și recomandări pentru dezvoltarea unei e-infrastructuri cu scopul de a menține accesibilitatea informației pe termen lung și utilitatea informației științifice digitale în cele mai importante depozite.

Pentru a realiza o infrastructură care să ofere acces pe termen lung, de multe ori se aplează la un sistem deschis de tipul Open Archival Information System (OAIS) Reference Model (ISO 14721). Termenul de acces pe termen lung definește o perioadă de timp suficient de lungă care să vizeze impactul produs de schimbarea de tehnologie, inclusiv suportul pentru o nouă media și noi formate de date.

Pașii parcurși în e-Știință a modificat dramatic procesul de cercetare. Ciclul de scriere și citire a publicațiilor în secolul trecut care utiliza un singur mediu pentru schimbul de informații a evoluat către o multitudine de resurse digitale. Aceste noi oportunități promovează cercetarea multidisciplinară și permit reutilizarea rapidă a

informațiilor. Pentru a avea însă acces la cercetări este necesar a întreprinde eforturi coerente și concrete pentru a prezerva înregistrările digitale.

5.1.1. Organizații implicate în prezervarea conținutului digital

Alianța Europeană pentru Accesul Permanent (European Alliance for Permanent Access) (<http://www.alliancepermanentaccess.eu/>) Alianța este o coaliție unică trans-sectorială a majorității deținătorilor de informație din știință. Printre membrii alianței se numără: European Science Foundation, Centre National d'Etudes Spatiales, Centre Informatique National de l'Enseignement Supérieur, Joint Information Systems Committee UK, the British Library și National Archives din Suedia.

Practic, membrii alianței sunt deja activi într-un număr relevant de proiecte de prezervare pe termen lung și acces cum ar fi CASPAR, DRIVER, PLASNETS, SHAMAN și DPE. Implicarea partenerilor comerciali în unele dintre aceste proiecte cum ar fi IBM, Microsoft și SUN este benefică și vor mări sfera de influență în domeniul alianței.

Membrii alianței sunt implicați deasemenea și în dezvoltarea unei noi generații a depozitelor digitale științifice, care iau în considerare eforturile europene de construire a unei e-infrastructuri cum ar fi GENESI-DR, EuroVO-AIDA, acoperind știința pământului și comunitățile în astrofizică. Totuși, prin definiție, proiectele au limite în ceea ce privește focalizarea pe un domeniu și timpul alocat. Alianța Europeană pentru Accesul Permanent își propune să furnizeze un cadru conceptual și strategic care să organizeze și consolideze eforturile individuale și să umple eventualele goluri dintre acestea. Ar trebui să se finalizeze o abordare unitară europeană și în același timp să se încerce o coroborare și cu alte proiecte din afara Europei.

De aceea alianța a adoptat un plan de lucru strategic pentru a ajuta consolidarea rolului său la nivel național European și internațional. Planul de lucru va stabili acțiunile care vor ajuta la generarea masei critice de utilizatori și va accelera progresul pe termen mediu (2008-2010).

Este general recunoscut că un răspuns la provocările prezervării digitale necesită abordări coerente pentru toate tipurile de informație codată digital. Diversitatea comunităților și a practicilor riscă o multiplicare a soluțiilor parțiale, favorizând apariția unor neconcordanțe.

Accesul permanent trebuie văzut ca fiind ca parte a unei activități de creare a unui depozit digital. Se impune o viziune coordonatoare asupra costurilor, beneficiilor, standardelor și tehnicilor care pot deriva de aici. Managementul de risk trebuie să se bazeze pe o abordare validă a viitorului.

5.1.2.Preocupări actuale pentru implementarea depozitelor digitale instituționale în universitățile din România

În cadrul unei universități, misiunea asigurării cu material documentar revine bibliotecii. Sistemul național de biblioteci din România cuprinde în cazul bibliotecilor universitare, biblioteci subordonate direct Ministerului Educației, Cercetării și Inovării, cele patru biblioteci centrale universitare din București, Iași, Cluj și Timișoara și celelalte biblioteci universitare subordonate senatelor universităților.

Dacă în cazul bibliotecilor centrale universitare situația financiară permite achiziții de material documentar și de baze de date, în cazul bibliotecilor universitare subordonate universităților situația nu este la fel de bună.

În practica achiziției partajate, a celorlalte țări, funcționează consorțiile. Prin aceste consorții se negociază prețuri accesibile la baze de date pentru toată comunitatea academică. Deși prețurile la publicațiile electronice cresc an de an, „criza jurnalelor” se simte acut în nevoile de documentare ale universităților. Nici până acum nu există în România un consorțiu al bibliotecilor universitare care să funcționeze. Consiliul Național al Cercetării Științifice din Învățământul Superior a achiziționat pentru întreaga comunitate academică românească baze de date importante pentru dezvoltarea cercetării științifice în anii 2007-2008.

În Universitatea Transilvania, biblioteca universitară oferă acces la trei baze de date: Springerlink, Mathscinet și Forest Science Database. Din anul 2004 biblioteca nu mai posedă abonament pe suport tradițional la nici o publicație străină, iar cele trei baze de date nu acoperă nevoile de informare din toate domeniile. O nemulțumire persistă în faptul ca aceste produse pot fi accesate numai de pe rețeaua universității. Majoritatea doresc ca accesul să se facă de acasă. Poate de aceea rezultatele cercetării calitative privind atitudinea cadrelor didactice arată că în proporție de 48,1% accesează bazele de

date de câteva ori pe an.

La nivelul Uniunii Europene, sistemul de editare științifică este văzut ca un element fundamental al sistemului European de cercetare. Referitor la noua paradigmă a comunicării academice și anume mișcarea „acces deschis” în România există puține inițiative.

La Universitatea Transilvania, publicațiile proprii: Buletinul Științific și alte trei reviste: Jurnal Medical Brașovean, Recent și PRO Ligno permit accesul la nivel de rezumat, dar nu sunt indexate în Google Scholar. Vizibilitatea și impactul cercetărilor cadrelor didactice este mic considerând indicatorii scientometrici calitativi.

Statisticile internaționale reflectă știința românească într-o lumină nefavorabilă. Producția științifică a României raportată anual, numărul de articole la 1 milion de locuitori este un indicator slab. Nici situația României în Europa Centrală nu este mai bună.

Universitatea Transilvania a făcut eforturi susținute și în anul 2004 este poziționată pe locul 13 în Carta Albă a Cercetării. Începerea construcției **Institutului de Cercetare-Dezvoltare-Inovare: Produse High-Tech pentru Dezvoltare Durabilă**, indică preocupări deosebite de dezvoltare a cercetării științifice.

Analiza impactului citărilor pentru profesorii universității nu indică performanțe remarcabile. Implementarea unui depozit digital ar putea fi soluția ideală pentru creșterea vizibilității producției științifice a universității.

Universitatea Transilvania ar putea fi pionierul depozitelor digitale instituționale din România.

5.2. Elaborarea modelului conceptual al sistemelor de informare și documentare cu conținut tehnic. Rezultate obținute în cadrul Universității "Transilvania" din Brașov

Pe baza informațiilor obținute în urma documentării asupra stadiului actual de dezvoltare a depozitelor digitale instituționale în lume, a modului lor de implementare și a cerințelor impuse de cercetarea cantitativă efectuată în cadrul comunității academice privind dezvoltarea unui depozit digital în Universitatea Transilvania, se va elabora o strategie de marketing al cărei obiectiv va fi creșterea vizibilității cercetării științifice românești.

5.3. Implementarea unui depozit instituțional digital la Universitatea TRANSILVANIA din Brașov

Membrii comunității academice din Universitatea Transilvania sunt preocupați de dezvoltarea cercetării științifice și participă la toate competițiile naționale și europene lansate pentru accesarea fondurilor destinate cercetării. Numărul de proiecte câștigate este mare, iar producția științifică a universității este remarcabilă. Acest mediu academic este unul în care, cu succes, se poate implementa un depozit digital.

5.3.1. Misiunea și obiectivele depozitului digital

Depozitul digital este un produs al tehnologiei informației pe care universitatea îl oferă membrilor comunității sale pentru gestionarea și difuzarea materialelor digitale create de către instituție și membrii comunității.

Acest produs este în esență un angajament organizatoric privind gestionarea materialelor produse de universitate, organizarea, accesul și distribuția lor dar și conservarea pe termen lung.

Prin acest produs se vor colecta și distribui toate documentele științifice create de membrii universității.

Obiectivul principal al depozitului digital este creșterea impactului cercetării științifice a universității prin promovarea producției științifice a universității prin depozitul digital instituțional și, în acest mod, vizibilitatea ei pe plan mondial și european și valorificarea pentru creșterea aportului românesc la dezvoltarea științifică și tehnologică.

Echipa de implementare a depozitului digital

Pentru asigurarea implicării persoanelor specializate și asigurarea unui parteneriat performant se propune ca următoarele departamente și direcții de cercetare să realizeze pentru început un **DEPOZIT DIGITAL PILOT**.

1. Departamentul de informatizare

Departamentul a fost implicat în multe proiecte universitare, dintre care proiectele de dezvoltare tehnologică, proiecte software (Transilvania):

- **AGSIS** – Aplicația de gestiune a studenților și informațiilor școlare.
- **Admitere**, Utilizarea aplicației de admitere în întreaga universitate, pentru toate

domeniile științifice și pentru toate formele de admitere a candidaților la nivel de licență, master și școală doctorală: admitere prin lucrări scrise, probe practice, teste grilă, concurs de dosare etc.

- **FRACS**, extindere aplicație pentru a corespunde noilor cerințe de raportare din universitate și de la Consiliul Național al Cercetării Științifice din Învățământul Superior a activității de cercetare științifică a tuturor membrilor comunității academice
- **Portal.unibv.ro**, pentru studenți: registratură electronică – gestiunea electronică online a cererilor și adevierințelor cu transferul automat în AGSIS; securitate bazată pe criptare și coduri de bare
- **Site-ul universității**
- Introducerea unui suport de tip **Single Sign On** pentru autentificarea și autorizarea accesului în rețeaua universității pentru cadre didactice și studenți, cu acces ulterior la toate aplicațiile dezvoltate – colaborare cu Departamentul de Infrastructură și Tehnologie care este în curs de realizare și a cărei finalizare este prevăzută în 2010.
- Introducerea **semnăturii electronice** pentru toate cadrele didactice din universitate cu finalizare în 2011.

Considerăm că acest departament este cel mai potrivit în a lansa această aplicație electronică. Departamentul are expertiza dovedită în implementarea platformei elearning utilizând Moodle (<http://moodle.org/>), care este un sistem de administrare a cursurilor (CMS – Course Management System), și are o licență de utilizare gratuită („open source”), cu peste 20000 instalări în 171 de țări, un excelent instrument destinat tutorilor, implicați în crearea cursurilor „online” și administrarea claselor de studenți.

2. Platforma de cercetare interdisciplinară ASPECKT, Platforma/Laborator de Analize Statistice și Previziune a fenomenelor Economico-sociale și Cercetări de marKeTing – **ASPECKT**. În cadrul acestei platforme integrate de cercetare științifică și dezvoltare aplicativă există condițiile obiective și subiective necesare pentru implementarea depozitului digital. Colectivul de cercetare implicat în această platformă va participa la implementarea DEPOZITULUI DIGITAL PILOT.

3. Departamentul de cercetare Produse mecanice de înaltă precizie și sisteme mecatronice, cu direcția Digitizare și tratament de imagine

Membrii acestui departament vor participa la implementarea DEPOZITULUI DIGITAL PILOT.

4. Departamentul de cercetare Legislație și proprietate intelectuală, cu direcția Legislație și dezvoltare durabilă.

Membrii acestui departament vor asigura consultanță cu privire la dreptul de proprietate a autorilor și dreptul de a arhiva documente în DEPOZITUL DIGITAL PILOT.



Fig. 5.1. Echipa de implementare a DEPOZITULUI DIGITAL PILOT, ASPECKT-Dspace.

Echipa de implementare (fig. 5.1.) va fi un parteneriat între:

- Facultatea de științe economice;
- Facultatea de inginerie mecanică;
- Facultatea de drept.

Criteriile de selecție a echipei de implementare sunt:

- Departamente prietenoase în raport cu misiunea DEPOZITULUI DIGITAL

PILOT;

- Diversitate în zonele disciplinelor;
- Diversitate de conținut sau formate;
- Exemple de gestiune a diferitelor forme de proprietate intelectuală;
- Arhivarea unor colecții de dimensiuni mici pentru început;
- Echipă cu legături puternice și încredere.

Membrii echipei au lucrat și în alte proiecte. Fiecare prezintă expertiză în activitățile desfășurate.

Planul de dezvoltare al produsului DEPOZIT DIGITAL PILOT

Planul va cuprinde următoarele acțiuni:

- Identificarea echipei de service;
- Identificarea primilor inițiatori care doresc să își arhiveze cercetările în depozitul digital;
- Identificarea unor colecții deja existente care să fie arhivate;
- Identificarea conținutului nou pentru arhivare;
- Identificarea conducătorilor;
- Dezvoltarea politicilor;
- Dezvoltarea structurii consultative;
- Identificarea personalului academic implicat în echipa proiectului;
- Identificarea personalului neacademic implicat în echipa proiectului.

Resursele umane pentru service și asistență vor forma două colective:

- Comitet de asistență;
- Comitet pentru web-design,
care vor avea drept ținte ale activității lor:
- Definirea colecțiilor ce vor fi arhivate;
- Definirea fluxurilor de lucru - în Dspace;
- Arhivarea colecțiilor deja existente;
- Sprijinirea utilizatorilor depozitului prin oferirea de informații telefonice sau (pe o linie) on-line;
- Crearea de pagini cu întrebările cele mai frecvente.

Producția științifică a universității cuprinde: dizertațiile masteranzilor, tezele de

doctorat, lucrările științifice elaborate de cadrele didactice (Figura 5.2.). Lucrările elaborate de cadrele didactice sunt lucrări publicate în volumele conferințelor, lucrări publicate în reviste, cărți (Figura.5.3.).

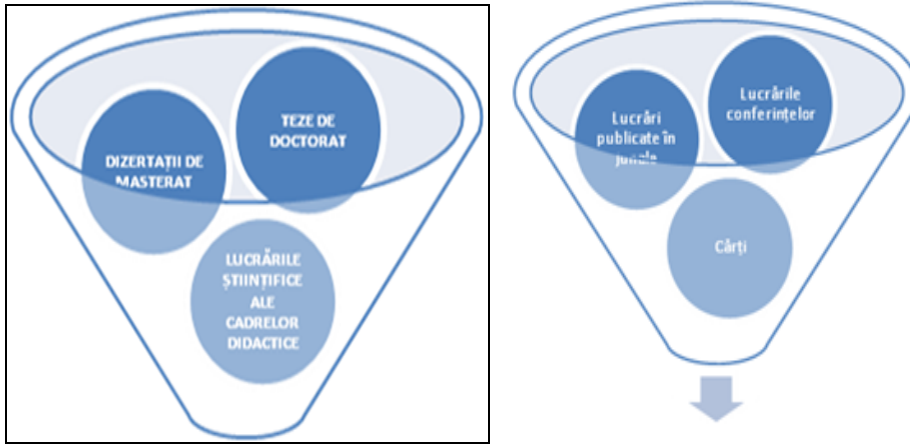


Figura 5.2: Producția științifică a universității. a Figura 5.3.: Lucrările publicate de cadrele didactice.

5.3.2. Tehnologiile pentru alegerea platformei și a softului

În urma unei analize tehnice se consideră ca Dspace este soluția sursei deschise pe care se va implementa DEPOZITUL PILOT în universitatea Transilvania.

Platforma de găzduire va fi platforma ASPECKT a Facultății de Științe Economice:

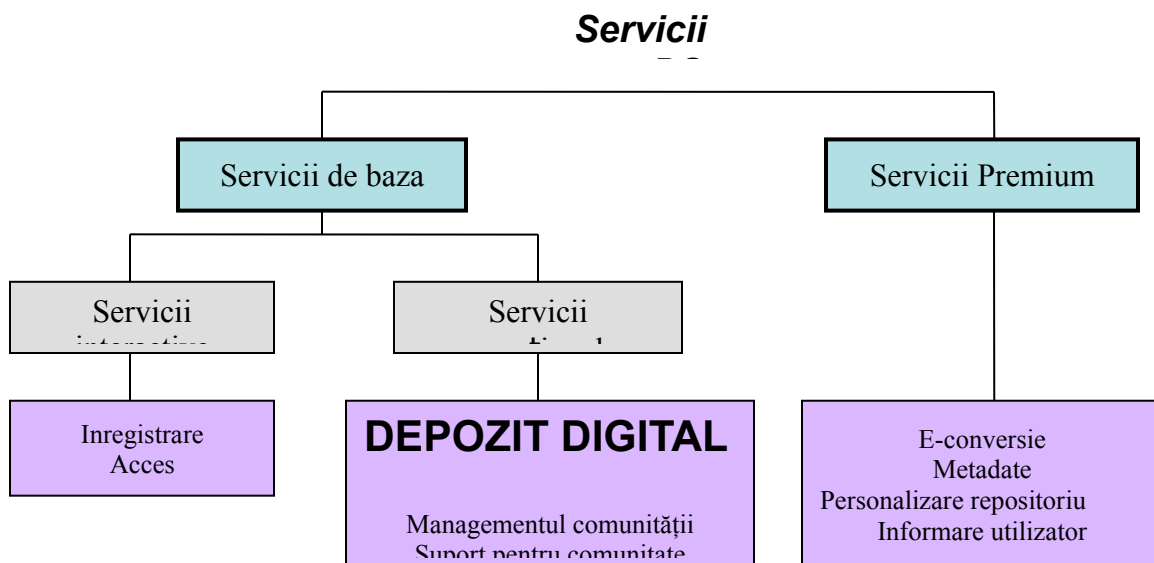


Figura 5.4. Serviciile oferite de platforma Dspace. Sursa (DSpace).

<http://aspekt.unitbv.ro/dspace/> în care DSpace a fost deja instalat.

Serviciile de bază oferite sunt:

- Managementul arhivării pentru a asigura o conservare pe termen lung,
- Stocare persistentă, inclusiv în caz de back-up și proceduri de recuperare,
- Atribuie fiecărui document un identificator unic, care să fie citat.

Dspace este softul tip sursă deschisă disponibil gratuit utilizat de cea mai mare parte a depozitelor instituționale. Oferă servicii de bază dar și servicii suplimentare premium, fig. 5.4.

Conținutul depozitului digital

Organizarea conținutului depozitului digital se poate face pe departamente de cercetare sau pe facultăți, la nivelul cărora se va lua și decizia asupra tipului de materiale ce vor fi depozitate.

Strategii de achiziție a lucrărilor ce vor fi arhivate în depozit sunt:

- inițierea echipei academice în auto-arhivare;
- redactare de instrucțiuni pentru ca autorii să-și arhiveze singuri articolele;
- redactare de instrucțiuni cu privire la managementul dreptului de autor, documentația de solicitare a permisiunii unui editor, în mod special, permisiunea în numele autorului;
- digitizarea copiilor pe suport tradițional este realizată de o echipă instruită;
- verificarea politicilor editoriale , a dreptului de autoarhivare.

Procedura de achiziție presupune următoarele etape:

- Autorii trimit pe email publicațiile dorite pentru arhivare.
- Echipa proiectului verifică dreptul de proprietate intelectuală și acordă dreptul de arhivare.
- Autorul își arhivează on-line publicațiile.

Dezvoltarea politicilor de achiziție, distribuție și menținere se face prin:

- Politicile de decizie se referă la conținutul depozitului. Trebuie stabilit ce documente se acceptă pentru a fi arhivate, cine poate arhiva documentele, cine va

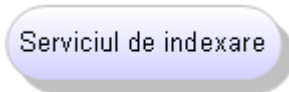
oferi metadatele. Se vor introduce și proceduri ce au în vedere calitatea lucrărilor: practic o comisie peer-review care se va activa în principal în cazul proiectelor de diplomă și care va avea obligația de a verifica și dacă proiectul nu este plagiat.

Managementul colecțiilor și a spațiului informațional se face ținând cont de serviciile oferite. Depozitul digital va oferi servicii de management al colecțiilor, de căutare, de indexare și stocare a documentelor:

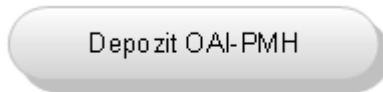
- La nivelul de prezentare avem serviciul interfeței cu utilizatorul și cel de publicare.



- Managementul spațiului informațional se face ca la orice bibliotecă digitală:

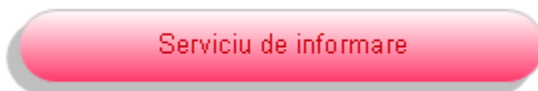


- La nivelul colectiv avem serviciul de culegere și agregare:



Serviciul de bază este acela de depozitare, stocare a documentelor.

Managementul sistemului funcționează cu un serviciu de autorizare și autentificare, serviciul de informare și serviciul de management:



Depozitul digital este diferențiat de alte colecții digitale prin următoarele caracteristici:

- conținutul este depozitat într-un depozit digital fie de către creator, deținător sau de terți;
- arhitectura depozitului destinează deopotrivă conținutul și metadatele;
- depozitul oferă un set minim de servicii de bază, de exemplu introducere, extragere, căutare, controlul accesului;
- depozitul trebuie să fie sustenabil și de încredere, bine susținut și bine gestionat.

Serviciile cheie pe care depozitele digitale trebuie să le acopere în câteva zone funcționale sunt (fig. 5.5):

- Acces crescut la resurse;
- Noi modalități de peer-review și de publicare;
- Managementul informațiilor la nivel corporatist (managementul înregistrărilor și sistemele de management ale conținutului);
- Partajarea datelor (reutilizarea datelor din cercetare, reutilizarea obiectelor didactice);
- Prezervarea resurselor digitale.

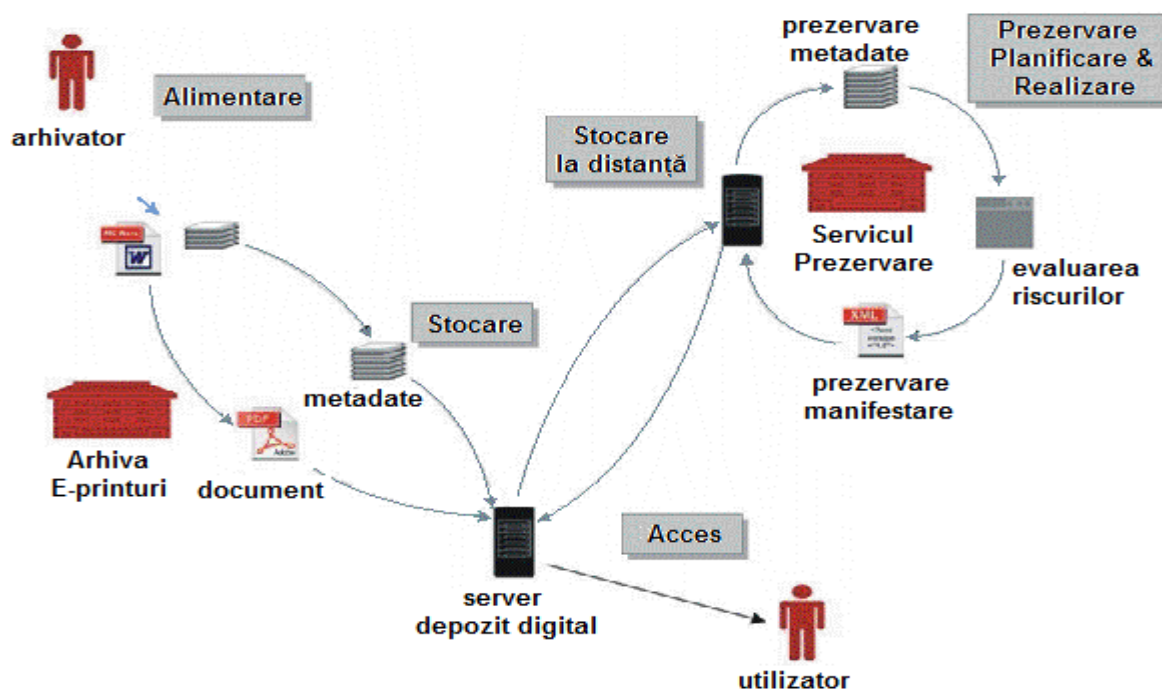


Fig. 5.5: Schema de funcționare a depozitului digital. Sursa (Kosson, 2009).

Depozitele digitale constituie de fapt un complex de servicii gestionate pentru a servi un scop și un anumit segment de utilizatori. Nu există o definiție care să stabilească distinct granițele a ceea ce este un depozit digital.

Lansarea depozitului digital

Produsul lansat este **DEPOZITUL DIGITAL PILOT** la Universitatea Transilvania - **ASPECKT-Dspace**. Lansarea acestui produs este concepută ca funcționând prin intermediul următoarelor etape:

- Trimiterea de materiale publicitare la toate departamentele de cercetare ;
- **Strategia de promovare a depozitului digital ;**
- **Strategia de profilare.** Această strategie este concepută pentru a convinge cercetătorii de beneficiile depozitului. Comunicarea informațiilor privind avantajele și beneficiile acestui serviciu sunt transmise utilizatorilor potențiali prin curier electronic al cărui conținut trebuie să evidențieze anumite elemente de bază, prezentate în exemplul de mai jos.

MODEL DE EMAIL TRANSMIS DEPARTAMENTELOR DE CERCETARE

- **Ce este ASPECKT-Dspace :** <http://aspekt.unitbv.ro/dspace/>
- Serviciul online dezvoltat arhivează materiale de cercetare publicate de membrii universității Transilvania;
- Aceste documente sunt disponibile de oriunde on-line și pentru oricine;
- Prin arhivarea cercetărilor în acest depozit, produsele științifice vor putea fi accesate pe o scară mai largă ducând la o mai mare vizibilitate și impact;
- Depozite similare cu ASPEKT-Dspace sunt dezvoltate în toată lumea;
- Documentele încărcate în Dspace sunt indexate în Google Scholar;
- Depozitul este un instrument care deschide accesul către cercetările din cadrul universității și le face disponibile pe scară largă și în întreaga lume;
- Nu este un mecanism de publicare sau susținut pentru reviste ci găzduiește materiale deja publicate;
- E destinat deschiderii literaturii științifice de cercetare în spiritul accesului liber la informație.

S-au stabilit comunitățile depozitului, adică tipurile de materiale ce se vor încărca în depozit și anume: Articole științifice, dizertații ale masteranzilor, lucrări studențești, proiecte de diplomă, rapoarte de cercetare ale doctoranzilor, teze de doctorat (Figura.5.6).

În Figura 5.6. se prezintă cazul căutării după autor.

Intr-o primă fază, se încarcă dizertațiile masteranzilor de la masterul de cercetare din domeniul marketing, și deci, autorii sunt masteranzii 2009. Valorificând acest exemplu, vor putea fi accesate, după același criteriu – acela al căutării după autor și alte categorii de materiale care au fost încărcate în aceeași bază de documentare și care se adresează unui anumit grup țintă vizat: lucrări studențești, proiecte de diplomă, rapoarte de cercetare.

În Figura 5.7. se prezintă lista completă a autorilor din categoria masteranzi la masterul de cercetare din domeniul marketing și ale căror lucrări pot fi accesate direct. Autorii sunt clasificați în mod ascendent sau descendent și, de asemenea, se poate opta pentru un anumit număr de rezultate ce sunt afișate pe o pagină.

The screenshot shows a web browser window with the URL <http://aspectt.unitbv.ro/dspace/>. The browser's address bar and tabs are visible. The page content includes a search bar with a 'Go' button, a 'Advanced Search' link, and a list of search criteria: Home, Communities & Collections, Issue Date, Author, Title, and Subject. Below this is a 'Sign on to:' section with links for 'Receive email updates', 'My DSpace authorized users', and 'Edit Profile'. A 'Help' and 'About DSpace' link are also present. The main content area features a 'Search' section with a text input field and a 'Go' button, and a 'Communities in DSpace' section with the instruction 'Choose a community to browse its collections.' and two links: 'Articole stiintifice' and 'Disertatii Masteranzi'.

Figura 5.6: Imagine din pagina web a depozitului digital ASPECKT-Dspace.

[DSpace at TRANSILVANIA University of Brasov >](#)

Browsing by Author

Jump to: [0-9](#) [A](#) [B](#) [C](#) [D](#) [E](#) [F](#) [G](#) [H](#) [I](#) [J](#) [K](#) [L](#) [M](#) [N](#) [O](#) [P](#) [Q](#) [R](#) [S](#) [T](#) [U](#) [V](#) [W](#) [X](#) [Y](#) [Z](#)
 or enter first few letters:

Order: Results/Page

Showing results 1 to 13 of 13

ANTOCI, Otilia
Badea, Anamaria-Irina
Badea, Anamaria Irina
BOITOR, Axenia Bianca
Boitor, Axenia Bianca
BUCUR BUCUR NICUSOR, Nicusor
BUCUR NICUSOR, Nicusor
CIRIC, Teodora
NAN, Elena
TIEREAN, Silviu
TIEREAN, SILVIU-HORIA
Mieracu Mieracu Timotei, Timotei

Figura 7: Structura depozitului pentru dizertațiile de masterat, căutare după autor.

În același mod se poate efectua căutarea și după titlul lucrării. În figura 5.8. este prezentat conținutul depozitului. Avem o imagine completă a documentelor, titlu, autor, data depozitării.

[DSpace at TRANSILVANIA University of Brasov >](#)

Browsing by Title

Jump to: [0-9](#) [A](#) [B](#) [C](#) [D](#) [E](#) [F](#) [G](#) [H](#) [I](#) [J](#) [K](#) [L](#) [M](#) [N](#) [O](#) [P](#) [Q](#) [R](#) [S](#) [T](#) [U](#) [V](#) [W](#) [X](#) [Y](#) [Z](#)
 or enter first few letters:

Sort by: In order: Results/Page Authors/Record:

Showing results 1 to 14 of 14

Issue Date	Title	Author(s)
Feb-2009	ACTIVITATE DE CERCETARE STIINTIFICA: MIXUL DE MARKETING IN DOMENIUL SOFTWARE - IT STUDIU DE CAZ: MICROSOFT VERSUS SUN MICROSYSTEMS	BOITOR, Axenia Bianca
Jul-2009	Activitatea de cercetare științifică II: Strategii ale mix-ului de marketing Studiu de caz: S.C. ALTEX ROMANIA S.R.L.	Boitor, Axenia Bianca
Jun-2009	ANALIZA DE MARKETING LA NIVELUL FIRMEI PLUS DISCOUNT ROMÂNIA	Badea, Anamaria Irina
Jun-2009	ANALIZA MEDIULUI INTERNATIONAL DE MARKETING AL COMPANIEI JAGUAR STRATEGII ALE MIXULUI DE MARKETING	TIEREAN, SILVIU-HORIA
Feb-2009	ANALIZA MIXULUI DE MARKETING BANCA COMERCIALA CARPATICA - BRD - GROUPE SOCIETE GENERALE	ANTOCI, Otilia
Feb-2009	ANALIZA STARTEGIILOR DE PIATĂ: COMPLEX BRAȘOV VS OCEAN FISH	TIEREAN, Silviu

Figura 5.8.: Structura depozitului după încărcarea dizertațiilor de masterat, căutare după

titlu.

Formatul prevăzut pentru încărcare este formatul .pdf deoarece el oferă mai multe avantaje, printre care cele mai importante ar fi cele legate de volumul mic al memoriei ocupate și de facilități foarte accesibile de securizare la copiere și imprimare, acces care ar putea rămâne la nivelul de decizie al autorilor.

Dspace odată instalat este foarte ușor de utilizat. În anexa 8 sunt prezentați toți pașii necesari pentru arhivarea documentelor.

Strategii de promovare a depozitului digital

Strategiile de promovare a depozitului digital se pot încadra în mai multe categorii: de atracție, de împingere și de consultare care pot fi implementate fie prin comunicare digitală – e-mail, internet etc., fie prin comunicare directă, bilaterală.

A. Strategia de atracție. Este strategia de a face depozitul atractiv pentru contribuabili potențiali. În acest scop un e-mail de antamare ar trebui să cuprindă mai multe elemente de teoria convingerii, ca în modelul de mai jos, ce pun în evidență avantajele promovate de depozitul digital.

MODEL DE EMAIL DE TRANSMIS MEMBRILOR COMUNITĂȚII pentru atragerea membrilor comunității să arhiveze lucrările științifice în depozit.

Beneficii pentru cercetătorii care își arhivează lucrările științifice în ASPEKT-Dspace:

- Cercetările sunt puse la dispoziția umanității pe o scară imposibil de realizat cu documentele tradiționale;
- Face documentele disponibile dispersându-le între mii de publicații academice dispersate în multe discipline;
- Documentele sunt indexate de motoarele Google și Yahoo;
- Cercetarea e stocată într-o centrală, spațiu de căutare și sunt distribuite într-o gamă largă de site-uri de cercetare;
- Depozitul reprezintă o alternativă sigură de a face documentele disponibile pe site-uri web personale;
- Accesul la arhive similare în toată lumea;
- Acces ușor pentru studenți.

B. Strategia de împingere. Această strategie are scopul de a crea o atitudine pozitivă față de depozit și de a arăta efectul pozitiv, odată ce materialele au fost arhivate. Mesajul trebuie să fie sugestiv, în acest sens și să conțină elemente destinate atingerii scopurilor strategiei de împingere, să pună în evidență ce este ASPECKT-Dspace, ca în exemplul următor.

MODEL DE EMAIL TRIMIS MEMBRILOR COMUNITĂȚII pentru a arăta efectul pozitiv al arhivării lucrărilor științifice în depozit:

- Serviciu on-line pentru a găzdui textul materialului de cercetare publicat de membrii universității Transilvania;
- Un produs destinat creșterii vizibilității și impactului cercetării științifice;
- Oferă condiții pentru cercetare liberă fără bariere de acces;
- Depozit în deplină siguranță pentru eternitate;
- Serviciu ce verifică acordurile dreptului de proprietate.

C. Strategia de consultare. În această etapă scopul este comunicarea bidirecțională. Acest lucru se realizează prin intermediul sondajelor, reuniuni, grupuri de lucru, de tip peer-to-peer de comunicare și alte forme de solicitare a feedback despre depozit.

Promovarea internă și marketingul depozitului digital

Promovarea internă și marketingul depozitului are în vedere atragerea utilizatorilor din cele două sensuri posibile: autori care furnizează materiale spre arhivare și utilizatori care se documentează pe baza materialelor conținute în depozit.

Această acțiune trebuie, în mod obligatoriu, să creeze anumiți vectori de comunicare și convingere a țintelor vizate prin oferirea de informații utile, clare, simple și cu mare impact referitoare la avantajele și beneficiile utilizării unui astfel de instrument de documentare și valorificare a cercetării științifice.

Acțiunile de marketing pentru promovarea depozitului digital, fără a fi exhaustive sunt:

1. Lansarea oficială a depozitului instituțional - DEPOZITULUI PILOT

2. Dezvoltarea la nivel de universitate, prin:

- Relaționare online: crearea unui blog pentru promovarea depozitului,

- Diseminarea experienței prin articole științifice, organizarea de workshopuri pentru promovarea experienței implementării depozitului,
- Organizare de cursuri de instruire.

3. Promovarea depozitului în lume

Pentru a crește vizibilitatea depozitului, acesta va fi înregistrat în cât mai multe liste de depozite și servicii existente:

- Repertoarul de Depozite Open Access (Directory of Open Access Repositories - DOAR);
- Registrul Depozitelor Open Access (Registry of Open Access Repositories - ROAR);
- OAIster;
- DRIVER.

Un lucru este sigur: depozitele instituționale vor reprezenta memoria unei societăți și un factor de progres în cercetare.

Sperăm că drumul deschis de Universitatea Transilvania să reprezinte un exemplu de bună practică pe viitor.

Concluzii privind implementarea unui depozit instituțional digital la Universitatea TRANSILVANIA din Brașov

- Membrii comunității academice din universitatea Transilvania sunt preocupați de dezvoltarea cercetării științifice și participă la toate competițiile naționale și europene lansate pentru accesarea fondurilor destinate cercetării.
- Numărul de proiecte câștigate este mare iar producția științifică a universității este remarcabilă. Acest mediu academic este unul în care cu succes se poate implementa un depozit digital.
- Depozitul digital este un produs al tehnologiei informației pe care universitatea îl oferă membrilor comunității sale pentru gestionarea și difuzarea materialelor digitale create de către instituție și membrii comunității.
- Acest produs este în esență un angajament organizatoric privind gestionarea materialelor produse de universitate, organizarea, accesul și distribuția lor dar și conservarea pe termen lung.

- Prin acest produs se vor colecta și distribui toate documentele științifice create de membrii universității.
- Obiectivul principal al depozitului digital este creșterea impactului cercetării științifice a universității prin promovarea producției științifice a universității prin depozitul digital instituțional.
- Realizarea unor depozite care să aparțină instituțiilor de cercetare, astfel ca munca cercetătorilor să fie vizibilă și crearea unui depozit central la nivel național nu solicită alocarea unor surse financiare considerabile, dacă ținem cont de costul unui sistem operațional, implementarea lui și echipa de dezvoltare.
- Fiind o tendință nouă în procesul de furnizare a rezultatelor cercetărilor academice, depozitele instituționale ar trebui să fie implementate în cadrul cât mai multor instituții, pentru ca astfel toți cei implicați în munca de cercetare să beneficieze de mijloace moderne de diseminare a cunoașterii.
- Includerea în OpenDOAR a primului depozit instituțional din România (fig. 5.9.)

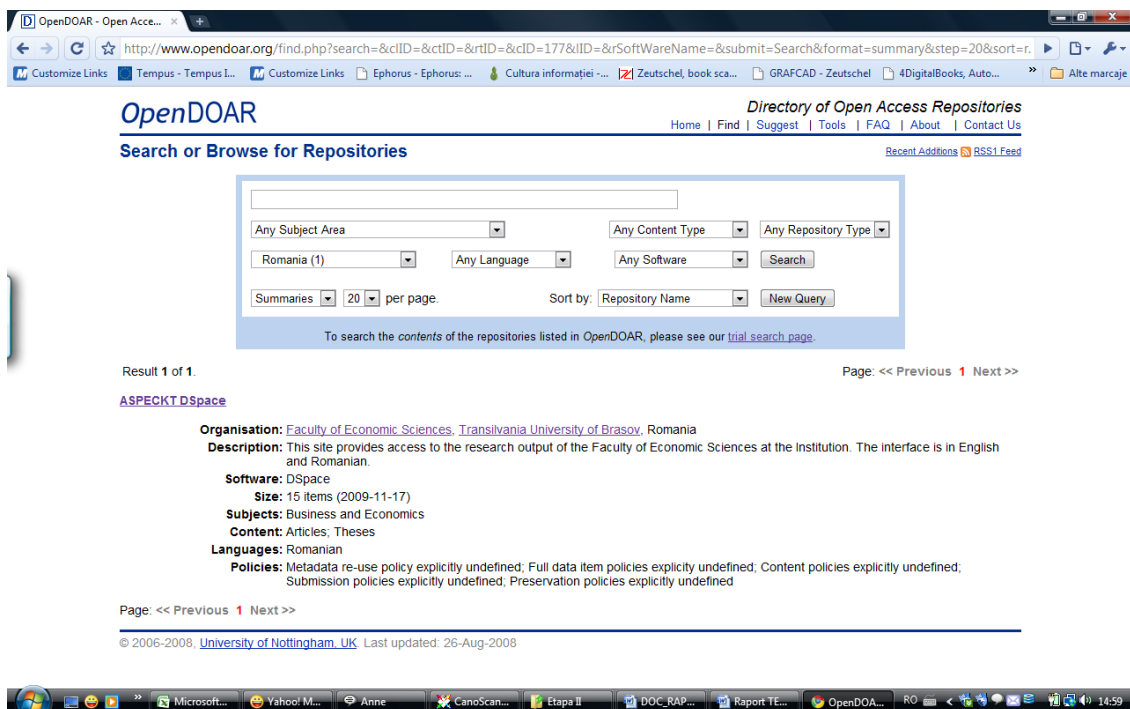


Fig. 5.9. OpenDOAR

6. Concluzii

Trăim într-o eră plină neprevăzut care aduce modele noi de dezvoltare a cercetării științifice, dar și dezavantaje pentru instituțiile academice care nu s-au implicat în proiecte de promovare a cercetărilor proprii, care riscă să intre într-un con de umbră în cazul în care nu își asigură o oarecare notorietate.

Datorită faptului că totul în prezent se schimbă cu „viteza luminii”, trebuie avute în vedere și nevoile de informare tot mai diverse ale comunităților academice. Acestea implică, printre altele și următoarele:

- folosirea surselor oferite prin intermediu Internetului (motoarele de căutare), care sunt mai ușor de utilizat și mai accesibile decât depozitele tradiționale ale bibliotecilor;
- realizarea unor depozite care să aparțină instituțiilor de cercetare, astfel ca munca cercetătorilor să fie vizibilă și crearea unui depozit central la nivel național
- sursele financiare, dacă ținem cont de costul unui sistem operațional, implementarea lui și echipa de dezvoltare.

Fiind o tendință nouă în procesul de furnizare a rezultatelor cercetărilor academice, depozitele instituționale ar trebui să fie implementate în cadrul câtor mai multe instituții, pentru ca astfel toți cei implicați în munca de cercetare să beneficieze de mijloace moderne de diseminare a cunoașterii.

Un lucru este sigur: depozitele instituționale vor reprezenta memoria unei societăți și un factor de progres în cercetare.

Sperăm că drumul deschis de Universitatea Transilvania să reprezinte un exemplu de bună practică pe viitor.

7. Note bibliografice si webografice

1. Banciu, D. coordinator, *Contributions to ICT progress - Digital Libraries from Concept to Practice*, Editura Tehnica, 2007, pp.: 77-94
2. Banciu, D., *Informatizarea structurilor infodocumentare*, Editura Ars Docendi, 2007
3. Bănică, L., *Metode și standarde la nivel de aplicație*, Ed.Universității din Pitești, 2005
4. Brașov, U. T. (2008-2012). *Plan strategic*. Brașov.
5. Coravu, R. (2008). Bibliotecile universitare și accesul deschis la informația științifică. (A. B. România, Ed.) *Revista română de biblioteconomie și știința informării*, 4 (3-4).
6. Campbell, J., *The Academic Library as a Virtual Destination*, EDUCAUSE Review, Volume 41, no.1, 2006, pp.: 16-31
7. Casarosa, V., *The European Digital Library - Discussion guidelines*, Workshop WS1, 2006, www.cimec.ro/ELAG/ELAG-2006-WS1.htm
8. Cristea, L., & Repanovici, A. (2007). Tehnologia informației-motorul trecerii la biblioteca digitală. În R. Angela (Ed.), *BIBLIO Brașov 2007. 2007*, pg. 317-318. Brașov: Editura universității Transilvania.
9. DSpace. (2009). *Dspace home federation*. Preluat pe septembrie 1, 2009, de pe DSpace Federation home page, <<http://www.dspace.org>>.
10. Gadd, E., Oppenheim, C., & Probetts, S. (Acceptat spre publicare). How academics expect to use open-access research papers. *Journal of librarianship and Information Science*
11. Gibbons, S. (2005). Faculty staff do not always understand the jargon used by librarians. The assurance of their work being found on the Web is better understood than the need for consistent metadata, and an unbreakable link better than a persistent URL. Preluat pe Iulie 29, 2009, de pe <http://docushare.lib.rochester.edu/docushare/dsweb/Get/Document-19725/amsterdam.ppt>

12. Hajjem, C. H. (2006). The Self-Archiving Impact Advantage: Quality Advantage or Quality Bias? University of Southampton, Department of Electronics and Computer Science.
13. Ignat, T. (2008). Biblioteca digitală și rolul său în activitățile de cercetare ale instituției pe care le deservește. (A. B. România, Ed.) *Revista română de biblioteconomie și știința informării*, 4 (3-4).
14. Ignat, T., & Repanovici, A. (2009). Institutional repository marketing: from research to knowledge transfer. *Advances in Marketing, Management and Finances, Proceedings of the 3rd International Conference in Management, Marketing and Finances, (MMF'09)* (pg. 107-110). Houston, USA: WSEAS.
15. Initiative, B. O. (2002). Preluat pe Mai 23, 2009, de pe <http://www.soros.org/openaccess/read.shtml>
16. Kevin Yank, *Build your own Database Driven Website using PHP & MySQL*, a guide published by SitePoint, SitePoint Tech Times newsletter, 2009
17. Kosson. (2009). *Depozite digitale*. Preluat pe septembrie 20, 2009, de pe Kosson comunitate virtuală: www.kosson.ro
18. OpenDOAR. (2008). *The Directory of Open Access Repositories*. Preluat pe Iulie 18, 2009, de pe <http://www.opendoar.org/>
19. Repanovici, A. (2009). Marketing Research about Attitudes, Difficulties and interest of academic Community about Institutional Repository, PLENARY LECTURE. *Advances in Marketing, Management and Finances, Proceedings of the 3rd International Conference in Management, Marketing and Finances, (MMF'09), Houston, USA, April 30-May 2, 2009, ISSN 1790-2769, ISBN 978-960-474-073-4, pag. 88-95* (pg. 88-95). Houston, USA: WSEAS.
20. Repanovici, A., & Turcanu, D. (2009). Qualitative and quantitative measures in marketing research for university library resource assesment. *Books of abstracts, Qualitative and quantitative methods in libraries QQML 2009*, (p. 73). Chania, Greece.
21. RoMEO, S. (2006). *RoMEO News*. Preluat pe Iulie 18, 2009, de pe <http://www.sherpa.ac.uk/romeo/>
22. Sherpa. (2006). *University of Nottingen*. Preluat pe Iulie 20, 2009, de pe <http://www.sherpa.ac.uk/>

23. SOROS. (2009). *Open access software*. Preluat pe Noiembrie 14, 2008, de pe <http://www.soros.org/openaccess/software>
24. Suber, P. (2005). Overview iin Open access- Unrestricted access to published research. *93rd Indian Science Congress* , 7-13.
25. Transilvania, U. (2009). *Institut de cercetare-dezvoltare-inovare: produse high-tech pentru dezvoltare durabilă*. Preluat pe septembrie 3, 2009, de pe http://www.unitbv.ro/institut_prodd/index.php?id=8
26. Transilvania, U. (2009). *Pagina web a universităţii Transilvania*. Preluat pe Septembrie 1, 2009, de pe Departamentul Informatizare: <http://www.unitbv.ro/Default.aspx?tabid=1222&language=en-US>
27. Ulman, L., *PHP for the World Wide Web*, Editura Teora, 2004
28. Welling, L., Thompson, L., *PHP and MySQL Web development*, 3rd Edition, Editura Teora, 2007
29. Witten, H., I., Bainbridge, D., *How to build a digital library*, Morgan Kaufmann Publishers, 2003
30. <http://allformp3.com/extension/what-is-the-dotx-files.html>
31. <http://www.altn.com/>
32. <http://answers.com/topic/doc-1>
33. <http://code.google.com/p/ocropus/>
34. <http://code.google.com/p/tesseract-ocr/>
35. <http://dot.extensionfile.net/>
36. <http://www.anydocsoftware.com/>
37. <http://en.wikipedia.org>
38. http://en.wikipedia.org/wiki/Rich_Text_Format
39. <http://euro.ubbcluj.ro/~alina/cursuri/birotica-practic/word/5-8.htm>
40. <http://euro.ubbcluj.ro/~alina/cursuri/birotica-practic/word/5-3.htm>
41. <http://ezinearticles.com/?Rich-Text-Format---Easy-to-Use-format&id=3204897>
42. <http://fileinfo.com/extension/doc>
43. <http://filext.com/file-extension/DOTX>
44. <http://interglacial.com/rtf/>
45. <http://kb.iu.edu/data/adnl.html>

46. <http://office.microsoft.com>
47. <http://office.microsoft.com/ro-ro/excel/HP100141031048.aspx>
48. <http://office.microsoft.com/ro-ro/help/HA100069351048.aspx>
49. <http://office.microsoft.com/ro-ro/products/HA101723691048.aspx>
50. <http://msdn.microsoft.com/>
51. <http://pumanet.codeplex.com/>
52. <http://ro.wikipedia.org/wiki/.doc>
53. http://ro.wikipedia.org/wiki/Microsoft_Excel
54. http://z.about.com/d/spreadsheets/1/0/a/2/-/-/excel_2007_screen_parts.gif
55. <http://wiki.mobileread.com/wiki/DOC>
56. http://www.adobe.com/devnet/flv/pdf/video_file_format_spec_v10.pdf
57. http://www.apple.com/downloads/macosx/productivity_tools/pdfocrx.html
58. <http://www.cvisiontech.com/>
59. <http://www.digitalsquad.ro/>
60. <http://www.driver-repository.eu/> DRIVER. (2009). *Networking European Scientific Repositories*. Preluat pe iulie 24, 2009
61. <http://www.docx.net/what-is-a-docx-file>
62. <http://www.docx.net/how-to-open-docx-files>
63. <http://www.docx.net/what-is-a-docx-converter>
64. <http://www.ecursuri.ro/cursuri/excel-prin-exemple.php>
65. www.edlproject.eu/conference/downloads/EDLconf_Gradmann.pdf
66. <http://www.encyclopedia.com/doc/1O11-RTF.html>
67. <http://www.eprintdriver.com/help/v5.0/COM/dllaux/DOCFmt.htm>
www.europeana.eu/ EUROPEANA project
68. <http://www.explaintechstuff.com/how-to-view-docx-files.html>
69. <http://www.file-extensions.org/dotx-file-extension>
70. <http://www.fileinfo.com/extension/docx>
71. <http://www.fileinfo.com/extension/dot>
72. http://www.jarte.com/help_new/document_file_formats.html
73. <http://www.htmlgoodies.com>
74. <http://www.microsoft.com/downloads/details.aspx?familyid=ac57de32-17f0-4b46-9e4e-467ef9bc5540&displaylang=en>

75. <http://www.musitek.com/X/default.html>
76. <http://www.scantron.com/cognition/>
77. http://www.scanstore.com/Document_Management_Solutions/
78. <http://www.scribd.com/doc/12831694/manual-microsoft-excel-complet>
79. <http://www.simpleocr.com/>
80. <http://www.techterms.com/definition/rtf>
81. <http://www.trainetrader.com/microsoft-excel-tutorial-series-the-basics-data-entry-and-navigation/>
82. www.theeuropeanlibrary.org/telplus -TELplus project
83. <http://www.trichview.com/features/files.html>
84. <http://www.w3.org>
85. <http://www.word.mvps.org/FAQs/Customization/CreateATemplatePart1.htm>
86. <http://www.word-tips.com/dotx-and-dotm.html>
87. <http://www.xml.com>