

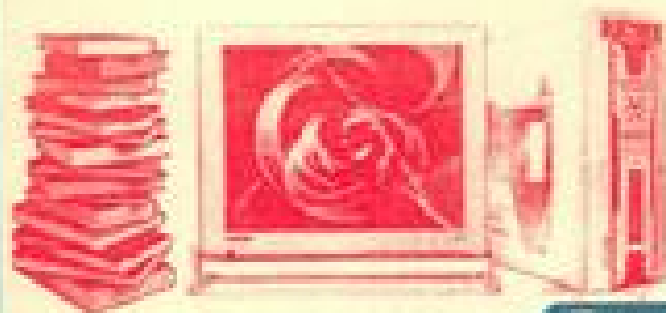


Полтавський університет  
своєвчотї кооператїї Українн

**ІНФОРМАЦІЙНІ РЕСУРСИ  
та ПОСЛУГИ:  
нові види, проблеми  
розвитку та  
використання**

*Матеріали міжнародної  
науково-практичної  
конференції*

8-9 червня 2004 року



Полтава  
РВВ ПУСКУ

руси и Республиканская научно-техническая библиотека смогли бы управлять использованием такой информации. РНТБ – патентной информацией, а ЦНБ НАН Беларуси научной литературой. При этом обязательным условием такой работы следует также считать организацию обучения потребителей информации, то есть открыть специализированный учебный центр.

## СТВОРЕННЯ ПОВНОТЕКСТОВИХ ДОКУМЕНТІВ ДЛЯ ЕЛЕКТРОННОЇ БІБЛІОТЕКИ

Влащенко Л.Г.;

Грищенко Т.Б.;

Нікітенко О.М., к.т.н.

Харківський національний університет радіоелектроніки

Одним з найважливіших показників зміни способу життя у ХХІ столітті буде розвиток та використання прогресивних інформаційних технологій у всіх сферах соціального життя та діяльності, рівень виробництва та споживання суспільством інформаційних продуктів та послуг. У цих умовах важливою справою є формування нової телекомунікаційної культури українського суспільства та вирішення проблеми якісно нової освіти.

Нові інформаційні технології призвели до принципів змін і в системі бібліотечного сервісу. Зникла необхідність у багатьох інформаційних формах бібліографічної діяльності й розвинулися зовсім нові форми пошуку та зберігання учбової та наукової інформації. На жаль, у більшості наукових бібліотек університетів мережеві комп'ютерні технології роблять тільки перші кроки.

Наразі все більшої популярності набувають колекції з повнотекстовими електронними виданнями, при цьому актуальною є проблема створення та доставки користувачу таких видань у електронному вигляді.

Тут під повнотекстовими документами мається на увазі будь-який документ у електронному вигляді (аудіо, графічний, текстовий).

### ***Вибір формату***

Графічні формати якщо й цікавлять нефаківців, то тільки з точки зору розмірів файлу.

З текстовою інформацією ситуація зовсім занедбана: добре розуміючи різницю між словом рукописним та друкованим, ми майже не друкуємо в мережі першоджерел у їх реальному форматі. Світові

деї й бібліотеки вже оцифровували більшість рукописів, які мають певну цінність, однак розміри отриманих файлів не дозволяють працювати з ними через Інтернет.

З грудня 1997 року у науковій бібліотеці Харківського національного університету радіоелектроніки (ХНУРЕ) було розпочато роботи зі створення своїх власних цифрових ресурсів. При цьому головними перевагами, що висувалися, були простота, зручність та ефективність технологічного процесу виготовлення електронних копій.

Ми бачимо, що вищезгадана проблема має розглядатися з урахуванням таких чинників:

- 1) швидкість виготовлення електронної копії;
- 2) розмір створеної копії;
- 3) зручність користування електронними копіями;
- 4) доставка електронної копії до користувача.

Розглянемо ці чинники детальніше.

Час виготовлення електронної копії поліграфічного видання складається з часу, що витрачається на сканування (вважаємо що перетворення твердої копії у м'яку відбувається за допомогою сканера) та обробки відсканованого матеріалу (обробки чернетки електронного документа). Час сканування залишається сталим і залежить тільки від типу використовуваного сканера. Отже, при виготовленні електронної копії суттєвим є час обробки відсканованої інформації.

Спираючись на чотирирічний досвід роботи зі сканування у бібліотеці ХНУРЕ електронні матеріали можна поділити на графічні формати запису *bmp, tif, eps, psx, djvu, gif, jpg*) та текстові (формати запису *html, doc, rtf, pdf*). Зрозуміло, що найменше часу витрачається на обробку графічних форматів *bmp, tif, eps, psx, gif, jpg*, бо тут після сканування треба тільки записати відсканований матеріал у відповідний файл. Трохи більше часу витрачається на створення файлів формату *djvu* через те, що процес обробки та запису у цьому форматі вимагає більше часу, ніж попередні. Найдовшою є обробка відсканованого матеріалу у текстових форматах, через те, що після сканування розпізнавання необхідно виправити помилки сканування.

Відомо, що зберігання відсканованого матеріалу потребує певних об'ємів пам'ятовуючих пристроїв, тому при створенні повнотекстових документів треба враховувати й цей чинник. Для з'ясування розмірів форматів було відскановано одну сторінку формату А4 з наукового журналу.

При цьому умови сканування були такі: формат А4, розрізняльна здатність *300 pdi*, режим сканування *line art*.

На рисунку 1 наведено діаграму розмірів відсканованої сторінки. З цього рисунка випливає, що за розмірами найбільш придатним форматом є *djvu*.

### Порівняльні характеристики форматів

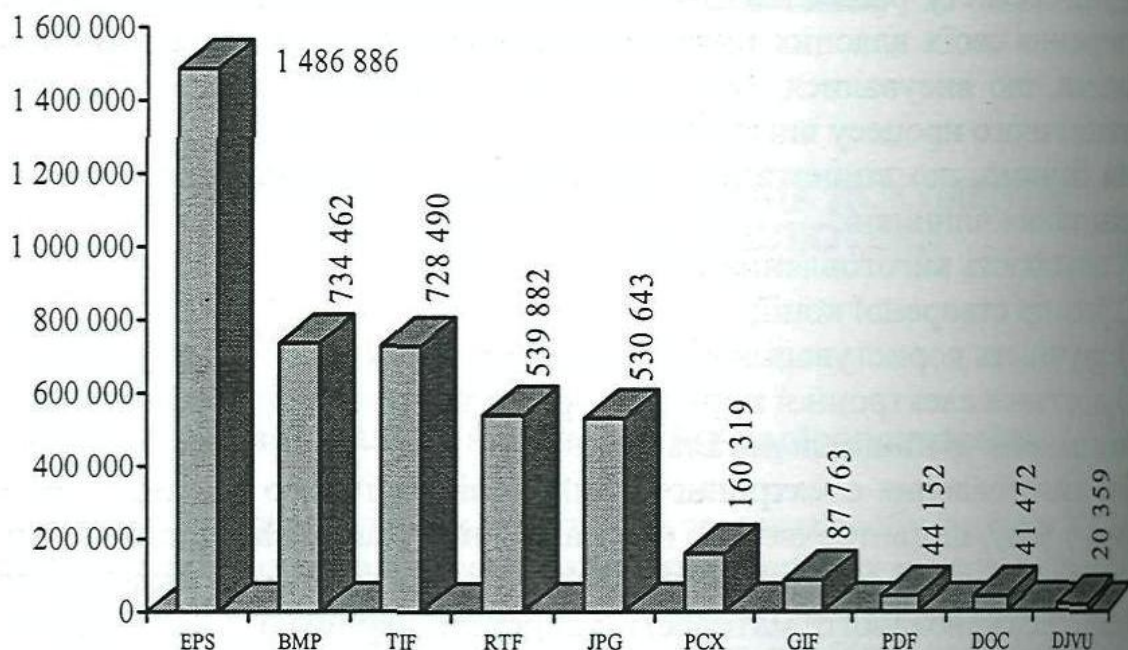


Рис. 1. Діаграма розмірів відсканованої сторінки

Однак усі графічні формати мають суттєвий недолік у порівнянні з текстовими: кожна сторінка відсканованого тексту міститься в окремому файлі, а це потребує додаткових зусиль для перегляду відсканованого документа.

Наразі здійснюються активні розробки нового графічного формату, який би задовольняв потреби мережі. Однією з перших закінчених розробок є формат *DjVu* фірми *AT@T*, який характеризується надзвичайно привабливими характеристиками – застосування нових ефективних алгоритмів стиску дозволило розробникам досягнути значного зниження об'ємів файлів. Фірмою розроблено програмні засоби для створення файлів формату *DjVu*, що є безкоштовними для некомерційного використання. Засобом перегляду таких файлів може бути звичайний *web*-браузер. Було здійснено роботи з освоєння технології створення документів формату *DjVu*.

### Огляд формату *DjVu*

Новий метод стиску графічних зображень, названий *DjVu*, створений спеціально для держання високоякісних копій відсканованих кольорових документів з високим ступенем стиску. Що дозволяє потім швидко передачу образу документа через мережу Інтернет в умо-

ліній з низькою швидкістю передачі даних (що особливо актуально для телекомунікаційних каналів країн СНД) відтворювати візуальну копію документа включаючи колір, шрифт, картинки і текстуру перу. Метод *DjVu* розподіляє зображення на текст і малюнок. Для тексту необхідний більш високий ступінь роздільної здатності, на відміну від малюнків та фону зображення. Малюнок же звичайно вищий за своїм складом, тому може бути закодований з вищим ступенем роздільної здатності. Потім застосовується новий код для максимізації ступеня стиску: дворівневий передній план зображення кодується за схемою *AT&T* у новий *JBIG2* факс-стандарт, для фону і малюнків використовується новий метод хвильового стиску *IW44*. Обидва методи використовують новий адаптивний арифметичний кодер, названий *Z*-кодером. Звичайна кольорова сторінка формату А4 із роздільною здатністю *300 dpi (dots per inch)* може бути стиснута до розмірів 40-60 Кб, що приблизно в 5-10 разів менше, ніж зображення у форматі *JPEG* при однакових рівнях ступеня стиску (рис. 1). Версія декодера зображення *DjVu* реалізована у вигляді програмного додатка для всіх популярних *web*-браузерів *Internet Explorer, Netscape, Opera*.

Стиск повнокольорової інформації про документ формату А4 до розміру середньої *web*-сторінки (46 Кб за даними на 1999 р.) теоретично-можливий. Враховуючи зростаючу суспільну потребу в доступі до «цифрових сигналів», здається дивним, що стандарт на графічні зображення цього призначення формується тільки сьогодні.

Формат *DjVu* – перший крок до «кольорового факсу» і його впроваджено на передачу, перегляд у мережі й роздруківку переважно текстових документів, для яких важливе значення має не тільки зміст, а й форма: колір та фактура пергаменту, відірваний куточок і сліди складання вчетверо, пляма після підпису й кругла пляма від винної пляшки поруч з печаткою. Архіви всього світу накопичили величезну кількість історичних паперів з неповторним колоритом такого чину.

Існуючі компактні формати *JPG, GIF*, факс-стандарт *CCITT* та *G* забезпечують достатній стиск, однак вузько спеціалізовані або фотографіях, або на чорно-білій графіці й тексті. Тому змішані зображення у їх виконанні виглядають таким, що важко прочитати. Розробники *DjVu* врахували негативний досвід створення «універсального солдата», їх розробка складається з трьох форматів «в одному лінійці». Розділ «обов'язків» всередині *DjVu* базується на простих зображеннях та фактах.

Текст та інші контрастні малюнки зручно читати при скануванні з роздільною здатністю не меншим за *300 dpi*.

Навпаки, невеличке розмиття фонові графіки навіть поліпшує сприйняття тексту. Тому фон без втрат для загального враження зберігається з розрізненням *100 dpi* в окремому шарі («*background*»).

Основна проблема – відокремити текст від фону, особливо якщо це кольоровий текст, і крім того, різнокольоровий. Здебільшого колір тексту в документах практично однаковий у межах одного знака. Це дозволяє зберігати кольорову інформацію про текст с розрізненням всього *25 dpi* (шар «*foreground*») (табл. 1).

**Таблиця 1. Розподіл у файлі формату *DjVu***

Шар	Пояснення	Розрізненість, <i>dpi</i>	Глибина кольору, <i>bits/pix</i>
<i>Mask</i>	Монохромна маска-трафарет	300	1
<i>Background</i>	Кольорове тло	100	24
<i>Foreground</i>	Кольори маски	25	24

Розподіл зображення на текст та тло (формування шару-маски) базується на так званій мультимасштабній кластеризації. Зображення розбивається на різнорозмірні вкладені сітки, в кожній комірці котрих відбувається розпізнавання текстових та фонових кольорів за максимальними пікам на гістограмі. Відокремивши текст від фону у найкрупнішій сітці, алгоритм переходить до уточнення на базі даних із сіток меншого розміру.

В *DjVu* для стиску фону, маски та кольорової інформації про маску застосовують різні алгоритми. Фон стискається вейвлет-алгоритмом *IW44* ( $4 \times 4$  *wavelets*), шар-маска, котра не містить кольорової інформації, стискається методом *JB2*, що є аналогічним до того, який застосовується у факсах. Кольорова інформація про текст також кодується *IW44*, але попередньо зменшується до *25 dpi*. Формат *IW44* є дуже близьким до нового стандарту *JPEG2000*, але менш вимогливим до системних ресурсів при декомпресії зображення під час перегляду.

Новий формат має багато застосувань: онлайнві книжкові магазини, картографічна інформація й навіть *e-хіромантия*, де надіслана за поштою фотографія долоні обробляється подібним чином.

Наразі цей формат вже здобув широке застосування у бібліотеках різного профілю, в релігійних організаціях, у радіоаматорських колах, використовуючи формат *DjVu*, видаються онлайнві математичні журнали.

Перехід до *DjVu* з його чітким текстом, на думку експертів, почнеться з сайтів ЗМІ, котрі копіюють свої паперові видання.

Формат *DjVu* дозволяє швидко переглянути матеріал у відкритому вигляді, й вже потім вирішити, чи варто його зберігати.

Якщо врахувати, що сторінка чорно-білої графіки з текстом має об'єм у форматі *DjVu* біля 30 Кб, а у кольорі біля 60 Кб, то стає зрозумілою економія часу та грошей.

Досить об'єктивна оцінка якості в порівнянні з уже відомими форматами показала, що незначне погіршення якості на кольорових зображеннях повністю компенсується ступенем стиску, котрий сягає тисяч разів, а на чорно-білих зображеннях практично не видно. Можливі конкуренти у вигляді *tiff, gif, jpg* сильно програють в обсязі.

До речі, популярний *jpg* зовсім непридатний для чорно-білих сканованих зображень.

— Саме на чорно-білих схемах і текстах, а поліграфічні видання більшого такими і є, перевага *djvu* є колосальною.

— Таким чином свій вибір ми зупинили на форматі *DjVu*, який пропонується американською компанією *LizardTech*. Формат *DjVu* дозволяє здійснити надто велике стискання зображень високої роздільності.

— Основні переваги цього формату:

1) — доступ до цифрової колекції по мережі Інтернет/Інтернет з використанням стандартного програмного забезпечення (необхідно лише встановити додатковий модуль для браузера);

2) — висока якість та малий об'єм зображень будь-яких видів (30 Кб для чорно-білого зображення формату А4 з роздільністю 300 dpi; 80-100 Кб для такого ж повнокольорового зображення);

3) — повне збереження виду видання;

4) — орієнтація на середовище Інтернет/Інтернет й простота забезпечення навігації всередині публікації.

Недоліком такого рішення є те, що сторінка у *DjVu*-форматі є зображенням і не дозволяє використовувати будь-який пошук. Власне, формат *DjVu* дозволяє зберігати текстову інформацію після процедури *OCR* (оптичного розпізнавання символів).

Зрозуміло, що бібліотеки мають обмаль коштів, яких ледве вистає на комплектування, тому бажано мати програмні пакети за мінімальною ціною. На щастя, в мережі Інтернет існують безкоштовні пакети, які дозволяють відтворювати повноцінні зображення у форматі *DjVu*, включаючи й багатосторінкові файли, перелік таких пакетів ведено в таблиці 2.

Таблиця 2. Програмні продукти для створення файлів формату *DjVu*

Назва	Обсяг, Кб	Призначення
<i>DjVuShop 2.0</i>	1613	Створення, редагування та перегляд односторінкових документів формату <i>DjVu</i>
<i>DjVuSolo 3.1</i>	2228	Створення, редагування та перегляд одно- й багатосторінкових документів формату <i>DjVu</i>

Назва	Обсяг, Кб	Призначення
<i>DjVuMulti 3.0</i>	2099	Утиліти командного рядка для створення багатосторінкових документів формату <i>DjVu</i>
<i>DjVuSDK2</i>	1907	Утиліта командного рядка для перетворення односторінкових документів формату <i>DjVu</i>
<i>DjVuWebBrowser</i>	1805	Перегляд документів формату <i>DjVu</i> безпосередньо в браузері

### Методика створення файлів у форматі *DjVu*

Маючи набір програмних продуктів, що наведено в таблиці 2, користувач може створювати файли формату *DjVu* будь-якої конфігурації.

Алгоритм створення файлу формату *DjVu* наведено на рисунку 2.

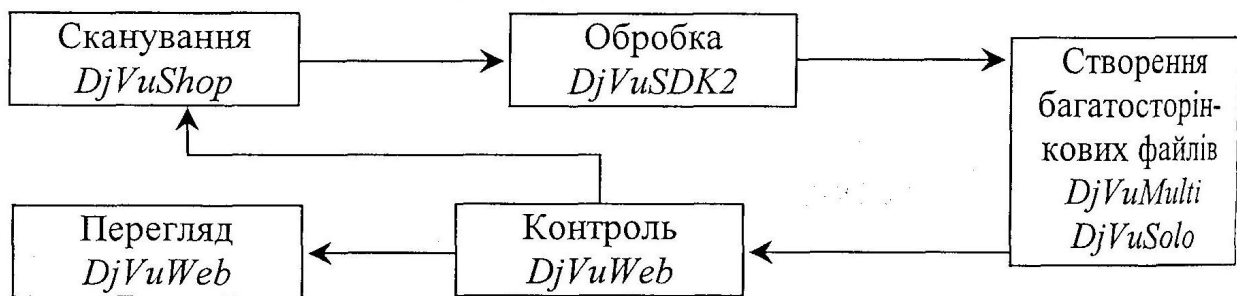


Рис. 2. Методика створення файлу формату *DjVu*

Сканування твердої копії та створення власне *DjVu*-файлу здійснюється за допомогою програми *DjVuShop 2.0*, яка зарекомендувала себе з кращого боку ніж програма *DjVuSolo 3.1* за часом створення файлу формату *DjVu* з використанням комп'ютерів *Pentium 120-150*.

За необхідності використовується обробка відсканованих файлів обертанням їх на необхідний кут. Ця операція здійснюється за допомогою утиліти командного рядка *DjVuSDK2*.

Після сканування та обробки всього необхідного матеріалу здійснюється створення багатосторінкових файлів. Це може бути як один багатосторінковий файл, котрий містить всі відскановані файли (операція *bundle*), так і створений додатково індексний файл (операція *join*). Обидві ці операції можна виконати, використовуючи або утиліти командного рядка *DjVuBundle* та *DjVuJoin* з пакету *DjVuMulti 3.0*, або програма *DjVuSolo 3.1*, яка дозволяє здійснювати ці операції у віконному режимі.

Через те, що створення багатосторінкового файлу – громіздка робота, необхідним є візуальний контроль створеного файлу, який здійснюється або програмою *DjVuSolo 3.1*, переглядом за допомогою будь-якого браузера, до котрого додано програму *DjVuWebBrowser*.



Після проведення всіх операцій створений багатосторінковий файл можна пропонувати для використання.

### **Висновки**

Виходячи з усього вищевикладеного в бібліотеці Харківського національного університету створюються повнотекстові електронні документи на базі формату *djvu*.

Але залишається ще одна проблема: доставка електронних документів читачеві.

Випробувані технологічні процеси створення файлів формату *DjVu* вели свою простоту, ефективність та економічність.

Для нас створення цифрових колекцій стало одним з важливих напрямків діяльності, яке ми збираємося розвивати й надалі.

Використання формату *DjVu* можна порадити тим, хто цінує час, прагне і якості. Не варто боятися нового формату, все було колись новим, краще просто спробувати самому, а потім вирішити, чи варта вина того, щоб її застосовувати.

## **БІБЛІОТЕКА ВИЩОГО НАВЧАЛЬНОГО ЗАКЛАДУ В СИСТЕМІ ДИСТАНЦІЙНОГО НАВЧАННЯ**

Щенко Т.Б.;

Щенко О.М., к.т.н.

Харківський національний університет радіоелектроніки

Постіндустріальне або інформаційне суспільство відрізняється виключно швидким розвитком інформаційних та комунікаційних технологій, можливості яких стають безпрецедентними для ефективного рішення багатьох професійних, економічних, соціальних та побутових проблем. Застосувати ці можливості можуть лише ті члени суспільства, котрі матимуть необхідну компетентність, що дозволяє адаптуватися в новому інформаційному просторі, зберігаючи свою конкурентність, використовувати переваги глобалізації. Мова йде про необхідність зміни змісту освіти на базі нових можливостей освітнього середовища, про засвоєння інформаційної культури, яку розуміють у вищій вияв людської освіченості, включаючи особистісні якості професійну компетентність.

Інформаційно грамотна людина визначається, як така, котра здатна усвідомити, коли їй потрібна інформація, яка вміє її знаходити, оцінювати та ефективно використовувати.