

Database backed Websites
D R A F T – Revision: 1.17

Thorbjørn Ravn Andersen

March 22, 2000

Dette dokument er under udarbejdelse.

Al tekst i dette skriftsnit er kommentarer til mig selv som jeg skal kigge på senere, og skal derfor vurdere
--

0.1 Abstract

This thesis discusses database backed websites; their current use as presentation engines, construction, pros and cons, list the current state of the art software, and how web sites in general function as *information presenters*. Hereafter it is discussed how end-users best can utilize such a system, without giving up their current software, allowing such websites to be used to ease the distribution of information between users, and let the web site provide *information sharing*.

as well as list my path to interest for meta-data and content-descriptive languages.

Technologies are discussed with an emphasis on Open Source software, and two sample programs - "Cactus" and "Consensus" are presented, and discussed. Cactus allows any user to do automated SSP on an Intranet without a human webmaster to convert documents and maintain links. Consensus allows an open group of people to maintain a set of documents on a webserver.

Stuff to place: Java has unicode support, multithreads.

Note: Framed text is unexpanded keywords which eventually will be replaced with prose.

Chapter 1

Overview

This section is rudimentary – it must be written at the end. Currently there is just keywords.

In this thesis I refer to the author Annie, and the layouter Larry.

1.1 What is a “database backed webserver”?

A web-server is a program which hands out files to any web-browser that asks.

A database “is a spread-sheet that several people can update simultaneously” (*greenspun*).

A database backed webserver then is simply a web-server which knows how to talk to a database server while servicing bypassing web-browsers.

1.2 Why use a webserver in combination with a database?

In one word: *Synergy*!

The web browser provide a simple, efficient and well known user interface available to almost every computer user in the world, and the database provide efficient dynamical access to data.

The combination of the two allow for the personalized Internet where a given webserver can give any user customized information upon request. Examples are:

- Newspapers - the user may specify that she wants in-depth coverage of international affairs, and do not want sports, and her personal web-edition will then just contain that. Payment may be per-article instead of per-paper.
- Maps/Airlines/Hotels - the user may request the best way to travel from A to B at a given time, see the route, and order accompanying reservations on-line.
- Real stores - the user may browse and order from the inventory. The ordered goods are then delivered by mail or similar.

- Virtual stores - the user may browse and order from a virtual inventory. The order is then dissected and forwarded to the appropriate suppliers to the virtual store.
- Banks - there is a growing need from customers to be able to do home banking. If the bank offer a browser based solution, they can support customers regardless of their choice of computer. The transactions done through such a system will most likely be stored in a database to avoid dataloss.
- Program development - The Mozilla project¹ has shown the benefit of up-to-date information regarding just about anything related to the source code. write about it

Since its inception in 1993 the web has grown to be the “good enough for most things” graphical user interface, capable of showing text and graphics as well as provide navigation. Due to the simplicity of HTML it is not hard to build a plain basic client² meaning that just about any modern computer with Internet capabilities have browsers available for them.

Provide background information: Quick overview of the technology [webserver, database, http, sql, dynamically generated content vs] with a graph. Give simple example. Explain *when* things happen (on-the-fly vs statically generated pages). Explain the advantages and disadvantages databases have over flat filesystems [transactions(Protection, atomicity), extra attributes, speed, indexed columns, complex queries, scalability(linear search in filesystems, multihost db's),]

1.3 The new role of the webserver

Compare original functionality with static HTML-files with content and a few CGI-scripts, to the current dynamically generated sites with many hits pr day. Since HTML is generated and provides little abstraction it is hard to use on a higher abstraction level, and with unreadable code in ASP and PHP3 it is harder. CSS was not the answer since it did not provide any abstraction from the presentation (seperate content from layout). Different needs of the user - wap/pdf/html/braille. CPU, storage, webmaster time prevents all formats being pregenerated.

Documents must be written in a format describing *content* and not layout.

¹...Mozilla project... – <http://www.mozilla.org>

²Plan basic client implies no Java, no Javascript, no cookies, no plugins, no cascading style sheets, no HTTP-1.1 and almost everything else that characterize the modern browsers from Netscape, Microsoft, Suns and Opera??. However, if a webpage is well-written it can still be displayed on such a plain basic client.

Chapter 2

Terms and concepts

““HTML is a SGML DTD””

[?] – [?]

This report uses a lot of abbreviations, many of which may not be widely known. This section lists most of them.

The list will be sorted

entity A named “unit” in XML/SGML which expands to a string. This is very similar to a macro without arguments in other languages.

FO Formatting objects. A generic description of the physical layout of a XML-document on paper. FOP converts this to PDF. *Several users have expressed doubt about this XSLT transformation*

SGML A family of languages. The *DTD* specifies which language it is. SGML can be formatted with FOSI or DSSSL style-sheets. “jade” can format to HTML, plain text and RTF

XML A family of languages. The *DTD* specifies which language it is. XML-documents can be transformed with *XSL*-style sheets to another XML-document (this process is called *XSLT*). XML is a subset of *SGML*

***ML** Common description of both XML and SGML, where the two have the same applications.

DTD The *Document Type Description* is the specification which describes the exact syntax of the *XML

XSLT Either the process of *transforming* a XML-document into another XML-document using a XSL-style-sheet, or the process of *formatting* a document into FO, which then can be further formatted in PDF

DOM *W3C’s initial model for representing an XML tree internally. Superseded by SAX (for java?)*

SAX Simple Application *interface for?* XML. An event based approach to XML-parsing and representation. Has an advantage in that the design allows processing to begin before the whole XML-tree has been read in. (*URL?*). Parsers and *processors?* which implement SAX can be selected freely, currently allowing for *20 different combinations of parsers and whashallicalem*.

HTML Hyper Text Markup Language. The “language of the web” which was originally designed by a physicist to present articles to other physicists. Was later hacked upon by Netscape and Microsoft to do things it was never originally meant to do.

DSSSL *SGML style sheets something*. Is used to render an SGML document to another format. “jade” can render to *FOT*, RTF, \TeX (must be post-processed with “jadetex”), SGML and XML. *Is this a standard?*

FOSI *Another style sheet for SGML*, which I do not know anything about. Apparently the US Navy uses it a lot. *Check in DocBook.*

SQL The Structured Query Language is the *standard* language for communicating with a database.

XSP XML Servlet Pages - A technique for combining code with XML in a single page, which is parsed and executed by the webserver when it is requested.

RTF *Rich Text Format* - a document description language *designed by microsoft*. Widely supported. May contain extensions to the original RTF-specification.

Jade A DSSSL-engine for SGML documents by James Clark.

JPEG An efficient method for representing photographs in computer files. Efficiency is achieved by discarding parts of the visual information which is hardest for the human eye to see.

GIF An image format well suited for computer generated images with few distinct colors, like icons. Is hampered by a patent on the compression scheme used and a maximum of 256 colors. Superseded by PNG.

PNG The successor to GIF. Is supported in newer browsers only.

CGI The Common Gateway Interface. The original way to generate pages and other files dynamically. A CGI-script can be written in any language supported by the web server.

IIS Microsoft Internet Information Server. The webserver from Microsoft.

ftp File Transfer Protocol. This is both the name of the *protocol* (the method) as well as the *program* implementing the protocol. ftp can transfer files between a client computer and a server computer, and is a standard on the Internet

http Hyper Text Transfer Protocol. This is the protocol a web browser needs to speak with a web server, in order to retrieve documents. Http is intentionally very simple.

TCP/IP This is the protocol which a computer needs to speak to connect to the Internet. It allows two computer to establish a data stream, where one machine pushes bytes in one end of the stream, and the other machine retrieves the bytes from the other end of the stream, without either computer being concerned about the underlying network.

SP A *what*?

Perl A scripting language which has become popular with Unix system administrators, and which “was there” when a need for scripting languages for CGI arose. Has eminent regular expressions.

ASP Microsofts version of a scripting language intertwined with HTML. Works on Microsoft web servers only.

JSP Java Server Pages. Suns version of Java intertwined with HTML. Requires a servlet-capable web server.

PHP The Internet version of a scripting language intertwined with HTML. Can run as a CGI script in most browsers, but as a module (which is faster) in at least Apache.

Apache Open Source webserver which runs on a lot of different platforms. Most Linux distributions include it.

Explain XML, XSL, XSLT, HTML (yes!), SQL, SAX, DOM, XSP and whatever fancy terms come to mind. Draw graph of what happens with an XML document.

Chapter 3

Dynamically generated documents

`{ 404 Document not found }`
`[?] - [?]`

!!! PUT SAMPLE OF FILENAME INSTEAD OF TT-FILENAME-/TT in SOMEWHERE.

3.1 The way it started: Handmade web-pages in a filesystem on a web-server

Back in the old days when the world wide web was designed, a web server could basically do only two things:

- Serve *static* files directly from an underlying file system. These files were either pages written in HTML, or images in GIF or JPEG formats.
- Call a program complying to the *Common Gateway Interface* with user-submitted parameters, and return its output as the document to send.

That was all!

Such web servers still exist - for example the NCSA server¹ (which has been unmaintained since 1996. NCSA recommend that the Apache server should be used instead). Even so the February 2000 Netcraft survey² reported that 17172 NCSA servers were active, and could be reached from Netcraft.

In August 1995 the number of webserver running the NSCA server was 10835³, meaning that number has increased slightly in that period. For comparison the total number of public webserver on the Internet in the same period grew from 19 thousand to 11 million, with Apache and Microsoft Internet Information Server accounting for 8.8 million of these.

Even with such a simple model, it works well for even *very* large web-sites which only generates few documents dynamically.

¹... the NCSA server... - <http://hoohoo.ncsa.uiuc.edu/>

²... the February 2000 Netcraft survey... - <http://www.netcraft.com/Survey/Reports/0002/>

³... webserver running the NSCA server was 10835... - <http://www.netcraft.com/Survey/Reports/9508/ALL/>

A well-tuned server which serves documents directly from the filesystem can reach impressive numbers. [Where was this - look for statistics](#). On a modern PC Apache has no problem saturating a 10Mbps Ethernet connection ([what about a 100Mbps connection](#)), meaning that even for very large and demanding sites the bottleneck will be determined by the hardware!⁴

On a related note: The ftp-server `ftp.cdrom.com` serves a lot of files every day. On March 14 2000 the transfer statistics⁵ said that the average output that day was 79.3 Mbit/s (with a peak of 95.8 Mbit/s). The average output on a yearly basis was 86.2 Mbit/s, which is 933 Gb/day in average). This single machine runs FreeBSD⁶ (see www.freebsd.org⁷ for details). Since an ftp-session have a notably larger overhead than an http-session, it is not unreasonable to expect similar performance from a suitably tuned web-server when serving static files.

CGI-scripts is another matter. For the occasional dynamically generated document CGI-scripts turned out to work well. All kinds of documents – images, webpages, progress reports – could be generated comparatively easy in the favorit language of the CGI-script author. Libraries to deal with the decoding of parameters, and encoding of the generated result were soon abundant, providing for many enhancements to the original “static web site” model.

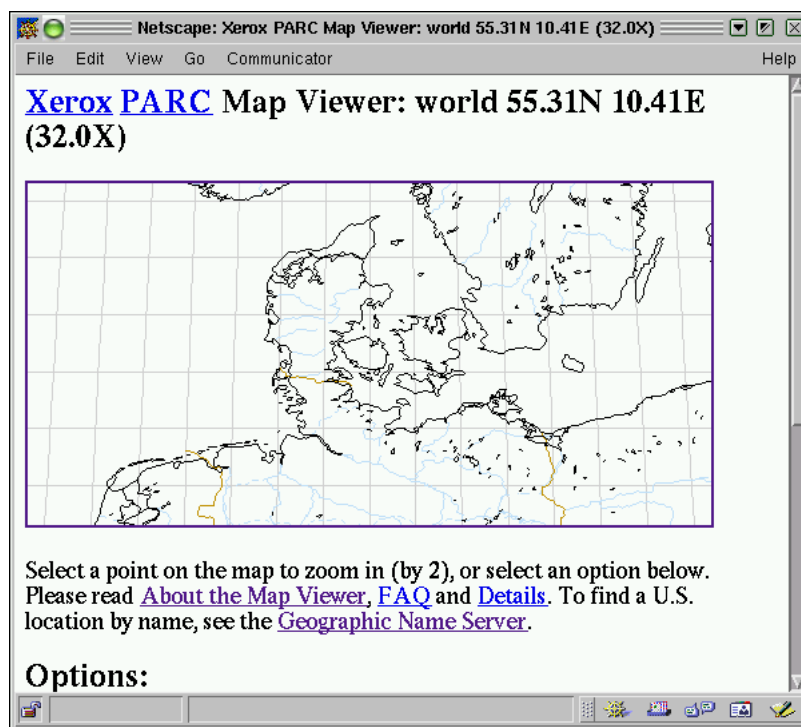


Figure 3.1: The original Xerox PARC Map Viewer

One of the first demonstrations of dynamically generated images was the Xerox PARC Map Viewer⁸ (see figure 3.1) which allowed the user to navigate a virtual atlas, where each “window” to the map

⁴ [Greenspun talks about SCHLOW webdatabaservers - write a bit about that](#)

⁵ ...the transfer statistics... – <http://www.emsphone.com/stats/cdrom.html>

⁶ ...This single machine runs FreeBSD... – <ftp://ftp.cdrom.com/.message>

⁷ ...www.freebsd.org... – <http://www.freebsd.org>

⁸ ...Xerox PARC Map Viewer... – <http://mapweb.parc.xerox.com/map>

was generated on demand. This initial quick hack has since been superceded by professional map companies which produce maps of an uneven quality on demand, like Yahoo Maps⁹ (USA only) and MapQuest¹⁰ (covers Europe too).

The problems with CGI-scripts is due to several reasons:

- **Script run as a subprocess** – the script is executed with the same permissions as the web-server itself, including access to `/etc/passwd` and other possibly sensitive files. Such scripts had to be *trusted* or - for ISP's - inspected before installation to ensure that the script would behave properly.

This has been addressed with the development of “safe languages” where the execution environment is guaranteed that the CGI-script is confined to a “sandbox”.

- **Incompatible platforms** – the CGI-programmer may not have access to a C-compiler which can generate binary code for the machine running the webserver.

This has caused interpreted languages like Perl to be very popular. Perl programs can be run immediately on a given platform without needing to be recompiled, and are for most purposes as fast as the equivalent programs written in C or C++. Recently Java Servlets (see section A.5 on page 99) have appeared as a popular and well-supported alternative to CGI-scripts, which is being supported by more and more web servers.

- **Slow** – the overhead just for invoking the CGI-program in a subprocess is substantial even for moderate load on the web server, since the web-server must “fork” a new process for each request.

This has been addressed by putting the execution environment inside the web server, using threads instead of processes, and caching compiled versions of scripts.

Several solutions to the above problems have emerged. Of these, Java Servlets are discussed in section A.5 on page 99, PHP3 scripts in section A.5 on page 99, and Perl (with the mod_perl acceleration module) in section A.5 on page 99).

Server side scripting never really took off before the PHP3 and ASP languages provided “persistent connections” to a database servers, which was unavailable to the standard CGI-scripts.

Persistent connections provide a real performance boost! Database connections are notoriously “expensive” to establish, so needing to open one for each and every CGI-script were a true performance bottleneck. Just by keeping the connection open (and for ASP-scripts - retain the particular database connection for the session) was extremely beneficial. That combined with a very simple way to switch between code and HTML plus that every user could use it without needing to ask the webmaster to install a potentially dangerous CGI-script, meant that the ability to generate webpages dynamically became available to everyone.

⁹...Yahoo Maps... – <http://maps.yahoo.com/py/maps.py>

¹⁰ ...MapQuest... – <http://www.mapquest.com>

3.2 When a site grows, it becomes hard to maintain

This section needs to be written.

Talk about the growth of the net, and broken links, both externally and internally
Jakob Nielsen. Start of next section should be rewritten

3.3 Navigational structure should not be maintained by the author

reference? Almost all dynamic web-sites which add and delete web pages run into the problem of maintaining site-integrity, regarding ensuring that all the “links” point to the correct document. For external document all you can do is to check occasionally that the page is still there, but for internal documents the webmaster has to do the updates manually. Updating HTML-documents manually is at best a tedious pain, because it is hard, repetitive, mindless labour, which should be left to a computer.



Figure 3.2: The Fyns GNU/Linux User Group web-site. The side bar and the top graphic was added automatically by a script by Tobias Bardino

Figure 3.2 shows a sample page from the Fyns GNU/Linux User Group web-site, showing the way FLUG chose to do it. The framework shown is added to the HTML-file with a small perl script which must be run to create each page.

This is a typical static framework for a site, where every web page basically is “inserted” in the home page since the set of links is identical for all pages. Usually this involves a link to a search

engine for the site.

In the original web servers this required that every web page was modified to include the snippet of HTML which produced this header, which is a relatively easy task with a modern scripting language like Perl. Note that care must be taken not to change the modification times of the files, since failing to do so invalidates local copies in browsers and webcaches, even though that the contents of the pages is unchanged. (Modern web servers allow a webpage to include other files with Server Side Includes¹¹, which takes care of all this - except doing it unconditionally on every page. The user must still remember to insert the appropriate command sequence in the server).

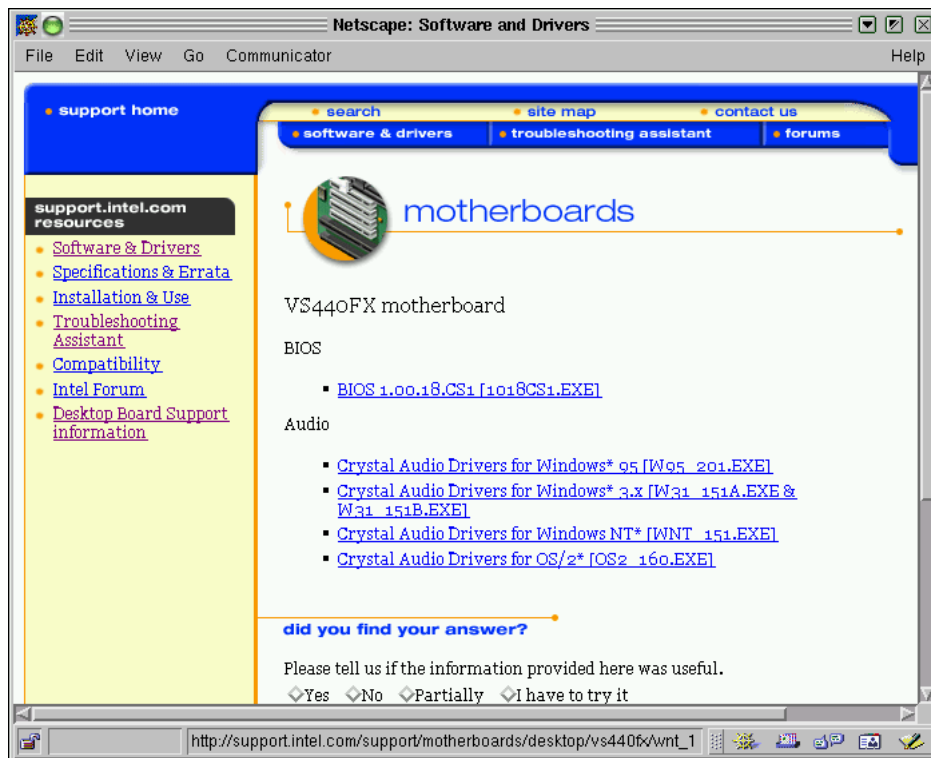


Figure 3.3: The web page for drivers for the Intel VS440FX motherboard

Many modern websites put a given page in a *context* where the navigational framework reflect this context. See figure 3.3 for a sample from Intel regarding the VS440FX motherboard¹², where the actual web page has a top bar with generic links, a column with links specific to motherboards in general, and a review form at the bottom asking whether this information was useful to the reader (this question is asked on every page with technical content). My personal experience with the Intel site is that their search engine is efficient, but that the link column is too uninformative - you often have to actually follow the link to see what is there. An expanded view would be nice.

This site layout require a bit more discipline for the web author writing each page, but is still easily implementable since it can still be implemented by including the appropriate HTML-snippet, as these are the same for all pages regarding this motherboard. A way to easily ensure this, would be to let the name of the physical file determine the virtual location in the web server hierachi - the

¹¹...Server Side Includes... - http://www.apache.org/docs/mod/mod_include.html

¹²...Intel regarding the VS440FX motherboard... - <http://support.intel.com/support/motherboards/desktop/VS440fx/softw>

shown file could have been named `motherboards/vs440fx/drivers.html`. This may be difficult to manage (if the webmasters work on different machines) and it is not very easy to change if the global layout of the web server changes.

In order to ensure integrity in the navigational framework for a site this size, it should be created automatically. Doing so requires meta-data about the individual pages in order to place them correct in the framework.

3.4 A good website needs meta-information about its documents

What to a webmaster is a nice, and well maintained web-site, is to a computer just a bunch of directories of files with bytes in them. In order to get any use of a computer in maintaining these, it is important to have easily accessible information about the desired functionality and the files to work with. Even though complicated rules can be constructed to extract information from webpages, the basic rule is still that

A human must enter the basic information for categorizing a given webdocument.

(Computers are simply not good enough yet to guess this themselves).

check spelling of altavista Incidentally this is also the reason why Internet search engines like AltaVista needs elaborate information extraction techniques in order to remain useful. AltaVista (which is discussed in section 5.9 on page 41) was the first Internet Search engine to cover *all* webpages. It was immediately a great success since the Internet had already grown to a size where help was essential to find any page if you didn't have a direct URL to it already.

Check precise way to do it and expand the text a bit. Web pages were considered in isolation or compared wi

It didn't take long for web authors (especially those with adult material) to realize that the best way to get their web pages frequently returned in a top position at AltaVista was by putting as many potential search terms in META-tags in each and every of their web pages, showing that just blindly extracting keywords *uncritically* is too simple a method. Additionally this drives users away – when they feel that they are not getting “the best results” from the search engine they will look for another one that can. The Google search engine (see 5.9 on page 42) was such an engine, and users quickly started to use it. Recently links to Amazon.com and Fatbrain.com have begun to crop up in the top of the search results on Google, and as a direct result users are going elsewhere. I have been recommended AllTheWeb¹³ as a good alternative.

AltaVista failed then, because they had no control over authors and the authors had a desire for “breaking the system”. For further reading, Douglas Hofstadter talks a lot about the impossibility of building an unbreakable system for automatic detection of bad input in [Hofstadter(1979)].

An example of a document generated from meta-data which has been automatically extracted, is the MIP “Recently Changed Pages” (see figure 3.4 on the next page) the original version of which I wrote while working for MIP. The script traverses three disjunkt set of web pages on the server - System pages, Sysop pages and User pages (one set per user) - and generates a list for each set sorted by title. Each file listed has its age in days next to it.

¹³...AllTheWeb... – <http://www.alltheweb.com>

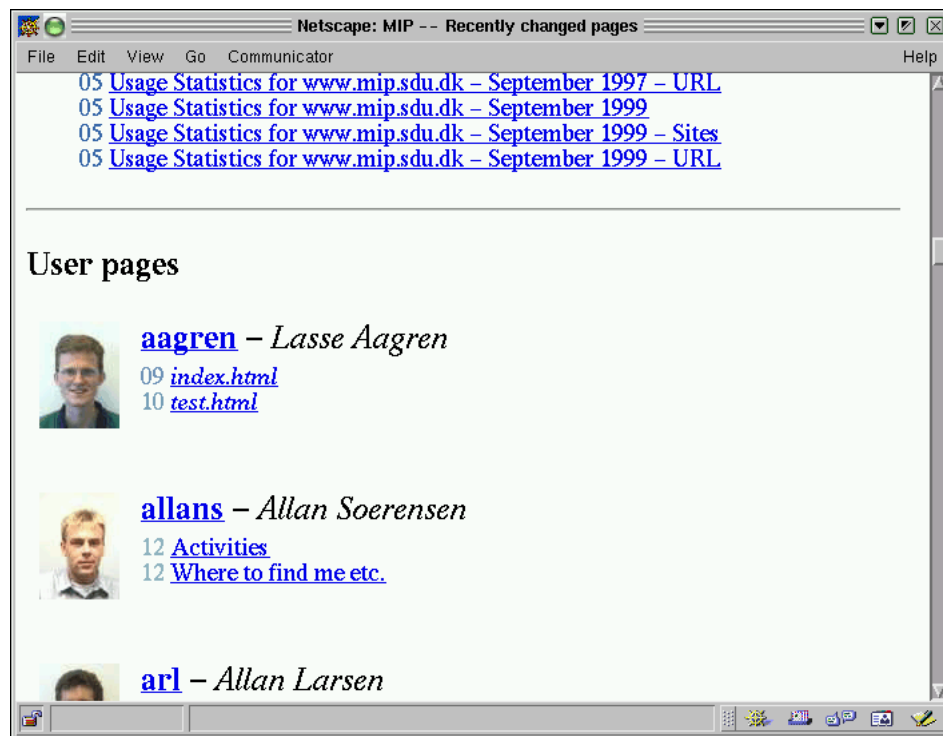


Figure 3.4: The MIP Recently Changed Pages. The last of the system pages and the first user pages are shown.

The meta-information extracted was:

- **Title** - used for the link, and sorting the entries in a set. Extracted from the content of each document, with a default of the filename if no title was present.
- **Age** - extracted from the underlying filesystem which registers the last change of the document
- **User** (for the user pages) - also extracted from the filesystem (*or was it from the URL?*). It is also used to look up the picture of the user.
- **Size** - *Did I use it? For anything? Check code*

To me this is just about all the useful information a flat Unix-filesystem can provide. Additionally nothing at all is guaranteed about the contents of a HTML-file - even the title is not even always there!

If more meta-data than the above is needed, the authors must be involved and as a part of their web-authoring, deliberately and carefully ensure that the meta-information needed by the automatic processes is correct and up-to-date.

Is there an Internet Library effort out there?

The Yggdrasil system (see section 9.1.3 on page 76) was my first attempt to generate a navigational framework from information extracted from webpages. Yggdrasil was to function as the automatic

webmaster on the intranet, in order to avoid having to assign staff to do so. Then the individual employee could publish information in form of a webpage, add a category either in the title or as a META-tag, and trigger the next update of the framework. This update would scan all web-pages, and extract tuples of (author, category code, publishing date, title, size) of those web-pages which had a category code, and generate a tree structure of web-pages to navigate the categories. Additionally lists of “Documents sorted by author” and “Documents sorted by date” were generated to help users locate documents.

This was on a closed Intranet, where the users were interested in using this as a tool. The incentive for “breaking the system” was very low, and the tool worked well.

Yggdrasil worked very well initially especially when its very low amount of meta-information is taken in consideration.



Figure 3.5: The British Yahoo site

An example of a good site built with meta-information is Yahoo¹⁴ (see figure 3.5 which started as a directory over web pages where the maintainers added meta-data to a lot of web-sites by categorizing them manually, and use this information to regularly generate static navigational pages. At a time Yahoo really suffered by a lot of broken links as the meta-data was not maintained, but today Yahoo usually links to a page that is there.

Internet Librarians? A comment?

¹⁴ ... Yahoo... – <http://www.yahoo.com>

3.5 Avoid chaos by keeping information together

Given there is a need for meta-data about the documents, the question is then *where* should the authors put the meta-data?

In the programming world, a very visible example of meta-data is the documentation for programs. Experience has shown that it is notoriously hard for programmers to keep the documentation up-to-date, since it is usually considered a part of the coding process that is not really necessary and definitively not felt to be a part of the “creative, fun” process of writing programs! If the writing of documentation was well established as an integrated part of program development, and the programming language itself helped as much as it possibly could, it would be easier for the programmers to keep the documentation synchronized with the code.

Experience has shown (references man month) that documentation should be as close to the thing it documents as possible. For programs, this means that the documentation should be *in the same file* as the code. Comments are usually just that – brief comments – and are often meant to explain *what* the code does instead of *why*? The abstraction level is not high enough.

grow terse - does not stand on its own - hard to maintain even for the author - abstraction layer: the idea

Andrew S. Tanenbaum writes excellent books about Computer Science, and his Operating Systems book [[Tanenbaum(1987)]] contains the complete source listing with line numbers of his Minix operating system, along with a cross-reference of all identifiers. That was the best the printing industry could do in 1987 (and is typical for several other printed versions of source code), and that was not good enough for teaching. It is very hard to get a grasp of what the code does, without long and careful studies, and much flapping back and forth. My guess is that today Tanenbaum would write a hyperlinked book readable in a browser, perhaps even using the Literate Programming techniques discussed in section ?? on page ??.

The great news is that the industry recognizes the usefulness of documentation on the web, and the need for programmers themselves to create such documentation. In order to gain wide acceptance such meta-data management systems must be *standards*, either as *de-facto standards* or defined by a standards body like ANSI or W3C.

The rest of this section detours into other areas of computer science, where having several views on the data in question has proven beneficial.

3.5.1 Javadoc - embedding web-information in programs

This is perhaps the most visible and generally available meta-data tool for Java programmers today. Javadoc¹⁵ is a tool that creates documentation in form of HTML pages from Java source code with embedded comments, where all definitions are parsed and hyperlinked to give a full overview of the Java source code in question. See figure 3.6 on the next page for a Javadoc rendered source code at ACME labs.¹⁶

¹⁵...Javadoc... – <http://java.sun.com/products/jdk/1.2/docs/tooldocs/solaris/javadoc.html>

¹⁶...ACME labs... – <http://www.acme.com>

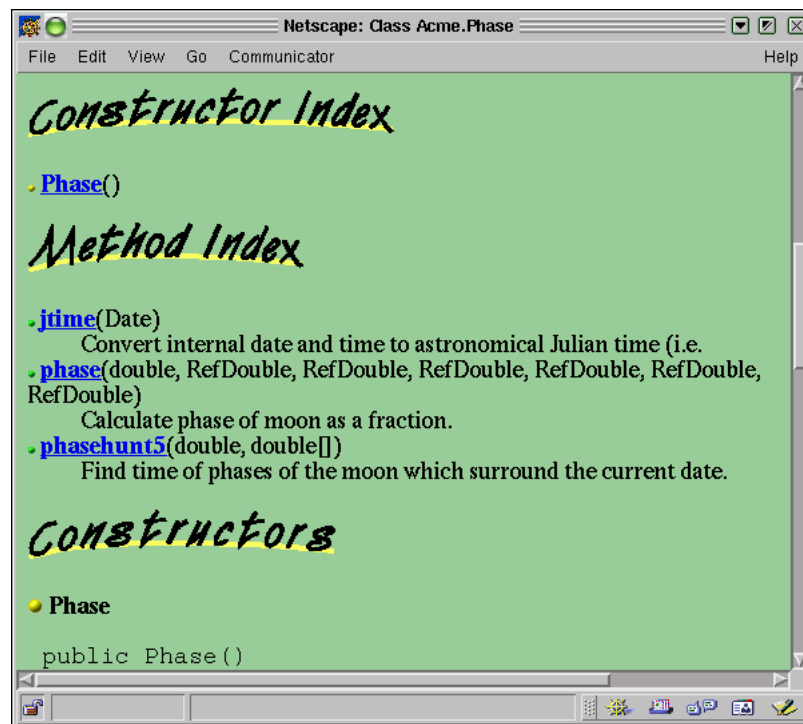


Figure 3.6: The ACME documentation of ??

JavaDoc is a specialized tool which is only suitable for generating web-pages from Java source code, but as Sun use JavaDoc for their reference documentation¹⁷ as well as ship it with every copy of Java Development Kit, it is a tool which is widely available. Programmers who need to document their code, will most likely use JavaDoc to do so.

Javadoc has also raised the expectations of the programmers, since they have become used to hyper-linked documentation in a browser to accompany any code they are to use. This tendency is a good step in the right direction.

Others recognize this too. On March 13, 2000 a an open letter pleading for the release of JavaDoc as OpenSource¹⁸ was posted to the Internet. They argued amongst other things that a number of bugs needed to be fixed, and new techology like XML/XSL should be employed too. Hopefully this will influence Sun to help the developers as much as possible by opening up this source too.

3.5.2 The Plain Old Documentation format for Perl

The perlpod¹⁹ format was designed to be simple to write and easy to use in Perl programs, and the overwhelming amount of documentation for Perl confirms that this goal was reached. This format allows you – with very elementary markup – to put documentation and code in the same file, and you may leave out either one.

¹⁷... use JavaDoc for their reference documentation... – <http://java.sun.com/products/jdk/1.2/docs/api/overview-summary>.

¹⁸...an open letter pleading for the release of JavaDoc as OpenSource... – <http://relativity.yi.org/WebSite/opensource-javadoc/>

¹⁹...perlpod... – <http://www.cpan.org/doc/manual/html/pod/perlpod.html>

All the info-files at Unixsnedkeren.dk²⁰ (the website for my small firm) have been written as POD files, and converted to HTML, with some finishing touches done with another Perlscript modified for the one I wrote for the FAQ for the Unix echo in the Danish Fidonet²¹. Code truly lives forever.

Due to the *ad-hoc* code in both “pod2html” and my own code, I am convinced that another approach to rendering POD-files should be taken, and I have submitted patches to the “pod2docbook” command in order to create DocBook XML in addition to the current DocBook SGML.

3.5.3 Literate Programming - Knuths approach to different views of the source

*“I believe that the time is ripe for significantly better documentation of programs, and that we can best achieve this by considering programs to be works of literature. Hence, my title: ”Literate Programming.” Let us change our traditional attitude to the construction of programs: Instead of imagining that our main task is to instruct a *computer* what to do, let us concentrate rather on explaining to *human beings* what we want a computer to do.”*

Donald Knuth – [Knuth(1992)] quoted from The Literate Programming web site²²

Donald E. Knuth has written the T_EX-system used to typeset this report, and documented it by publishing the source code to the entire system in four books. The T_EX-system has impressed since it is of an extremely high quality, both in code and documentation, and actually offer money to those who find errors in his code²³.

references to the T_EX book, MetaFont book, Literate Programming Book, nuweb (my observations)

In order to write both documentation and code of such high quality, Knuth developed his own method of coding *which he named “Literate Programming”*, in which the author works with a single file containing the documentation *as it is to be presented to the reader* written in a T_EX-dialect, with the actual code (with a few characters escaped) listed in named chunks along with their documentation, as it fits the author to present them. This file format is the web-format!

- The `weave` tool converts the web-file to a `tex`-file to be processed by T_EX, and printed.
- The `tangle` tool generates the actual code to be compiled, by joining chunks with the same name, and inserting them in other chunks that reference them (a simple macro expansion), eventually producing one or more flat files which can be compiled.

We definitively needs a sample of this!

Reference to literateprogramming.org

The problem with Knuth is that things that work well for him, does not work quite as well for the rest of us. No editors - not even Emacs - can help where a given file is separated in a lot of regions, being either T_EX-code or source code, meaning that the editor cannot provide the supporting functions which a programmer might have become accustomed to.

²⁰...Unixsnedkeren.dk... – <http://www.unixsnedkeren.dk>

²¹...the FAQ for the Unix echo in the Danish Fidonet... – http://www.fido.dk/faq/unix-faq/unix_r23.htm

²²...The Literate Programming web site... – <http://www.literateprogramming.com/>

²³...actually offer money to those who find errors in his code... – <http://truetex.com/knuthchk.htm>

You cannot code verbatim - some characters are reserved for other purposes and must be written differently, which may be very annoying to learn. Additionally the concept of chunks all over the document may be counter-productive if these are hard to navigate.

What else

Many tried this- cloned the functionality -etc - I tried it, and wrote a few programs with it.

Many people have experimented with the possibilities of a weave/tangle pair resulting in several software packages, notably:

FunnelWeb²⁴,noweb²⁵

<http://w3.pppl.gov/krommes/fweb.html#SEC3>

There is a webring for Literate Programming²⁶, which is an excellent place to look for further information. The Collection of Computer Science Bibliographies provides a search engine in literate programming publications²⁷.

Oasis has a page on Literate Programing with XML and SGML²⁸ *what about it? - note the references for working*

My previous experience with Literate Programs can be summarized as:

- The documentation gets bigger and better, simply because the printed documents look better that way. What would pass as a single line comment in a source file, looks almost pathetic when typeset. The full power of T_EX also encourages usages of illustrations and graphs.
- The program development gets cumbersome. You may have trouble using your favorite tools for editing, compiling and debugging. Users of most integrated development environments cannot use this model since the IDE does not provide hooks to provide this processing.
- A critical point is whether the intermediate files are visible to the author! In order to be usable, *all* derived files *must* refer to the original document in a transparent fashion, meaning that the user should not have to worry about intermediate files. (In the same way that the C-processor works under Unix).

My conclusion was that the Literate Programming paradigm does not pay off as a individual programmer, but may work very well for a larger programming team where good, current documentation is critical.

3.5.4 Rational Rose - the other way around

A software product which has been very successful in recent years, is the Rational Rose visual modeling tool²⁹ which is used by several people at MIP to do software development.

²⁴ ...FunnelWeb... - <http://www.ross.net/funnelweb/>

²⁵ ...noweb... - <http://www.eecs.harvard.edu/nr/noweb/>

²⁶ ...a webring for Literate Programming... - <http://www.webring.org/cgi-bin/webring?ring=litprog;list>

²⁷ ...a search engine in literate programming publications... - <http://liinwww.ira.uka.de/searchbib/SE/litprog>

²⁸ ...a page on Literate Programing with XML and SGML... - <http://www.oasis-open.org/cover/xmlLitProg.html>

²⁹ ...Rational Rose visual modeling tool... - <http://www.rational.com/products/rose/index.jtmpl>

Rational Rose provides for a lot of things relevant to a large software project like reverse engineering of existing code, but also for several views of the document³⁰ depending on what abstraction level the authors are working on. In this way they are able to provide the author an environment where the software model can be designed separately from the writing of the code.

Bo Nørregård Jørgensen told me that he would expect it to be beneficial to use Rational Rose in a software project when it has an underlying model and has 7 modules or more.

3.5.5 Compiling source code into an executable

Most software projects known to me does not consist of a single huge source file which would take ages to compile, but of several smaller files, which can be compiled individually and linked to an executable file. If a file with source code is changed, it is not necessary to recompile each and every file but only those which depend on this particular source file, and then relink the executable to incorporate the changes.

This is possible because the authors has provided meta-data about the system, regarding which source files the system contains, which commands to call to compile and link the source files, which libraries should be included, etc. Normally an Integrated Development Environment (like Borland C++ Builder³¹, Microsoft Visual C++³², or the Visual Age products from IBM³³) knows about these things, or a Makefile³⁴ is constructed to utilize the known relations built into “make” as well as allow the author to add new ones.

The real advantage comes in environments (multi-threaded, multi-CPU) where several compilations can run in parallel. This can be done safely by using meta-data to localize compilations which are independent from one another and execute these simultaneously. MIP have had a 24-CPU Silicon Graphics machine in which using a parallelizing “make” could reduce compilation times with a factor of 20. Not all software projects were easily parallelized - if the Makefile did not list all dependencies explicitly but expected that the normal compilation order would produce the unlisted files, the parallelized compilation would fail.

In effect, adding meta-data to the project allow the program development process to go faster.

3.6 Requiring multiple views of a document

Some web-sites have so many advertisements and auxiliary links in their layout, that they need an additional version which is well-suited for printing (i.e. only has a single ad). See figure 3.7 on the next page for an example from ZDNet³⁵, and figure 3.8 on page 23 for the “Printer-friendly version”.

³⁰ ...several views of the document... - http://www.rational.com/sitewide/support/whitepapers/dynamic.jtmpl?doc_key=35

³¹ ...Borland C++ Builder... - <http://www.borland.com/bcppbuilder/>

³² ...Microsoft Visual C++... - <http://msdn.microsoft.com/visualc/>

³³ ...the Visual Age products from IBM... - <http://www-4.ibm.com/software/ad/>

³⁴ ...Makefile... - <http://www.eng.hawaii.edu/Tutor/Make/>

³⁵ ...for an example from ZDNet... - <http://www.zdnet.com/zdnn/stories/news/0,4586,2468874,00.html?chkpt=zdhpnnews01>



Figure 3.7: The ZDNet presentation of a story – notice how little of the initial screen that is dedicated to the story about “Microsoft may try to boost WinCE, Linux-style”

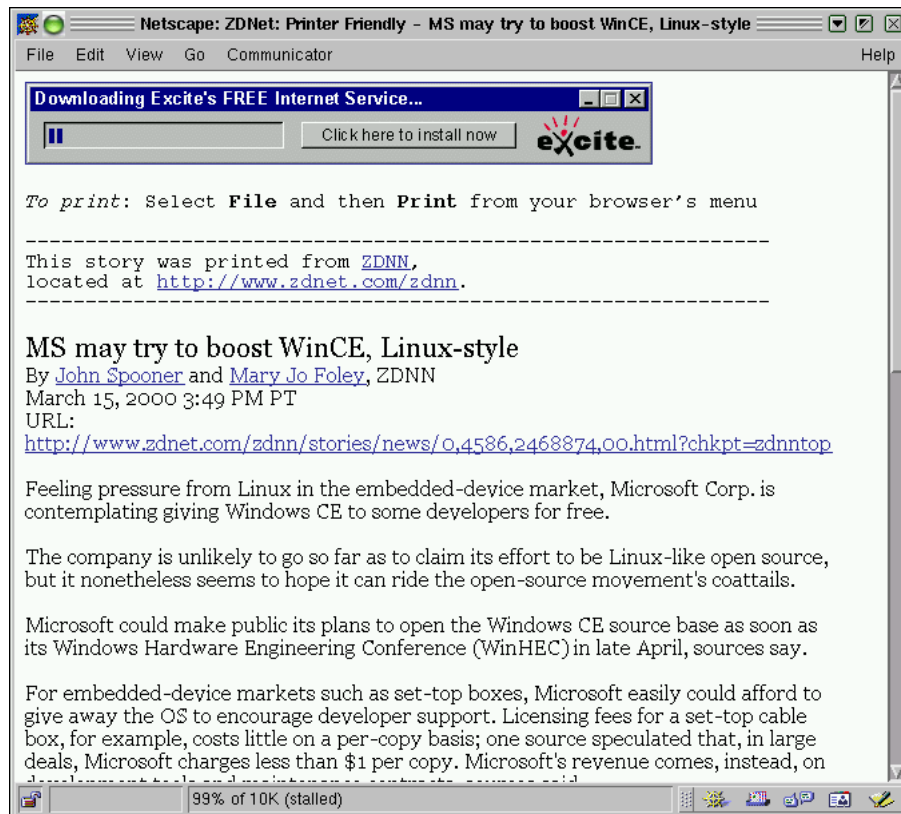


Figure 3.8: The "Printer-friendly version" of figure 3.7 on the facing page

It is interesting to notice that websites realize that they cannot just offer an advertisement ladden service, without allowing for a reasonable paper version. Since their HTML cannot “remove” the ads for printing, they have to offer two different version of each article.

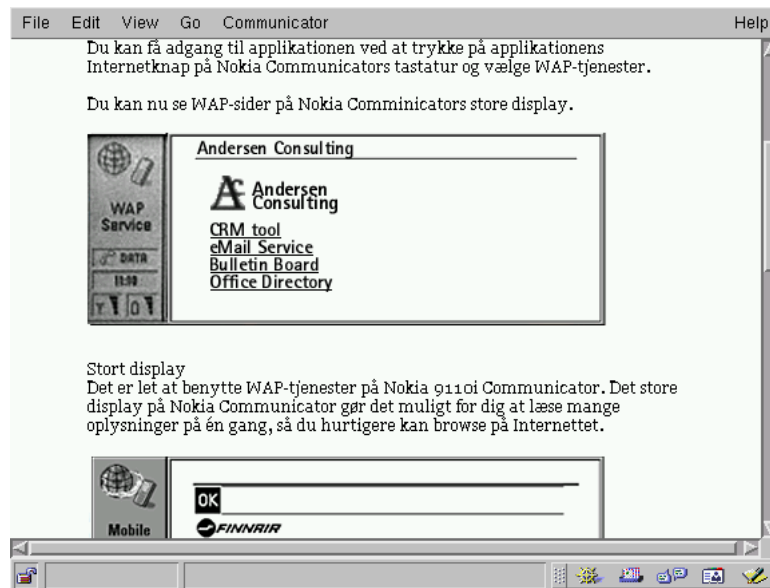


Figure 3.9: Browsing the Internet on the large [sic] display of a Nokia 9110

These are two different *views* of the same basic document. ZDNet are expecting their visitors to either see their documents on a modern display unit with high resolution and millions of colors rendered by a reasonably capable browser or print them out on paper. That will probably change within a few years when it becomes common for other devices than just traditional computers to be connected to the Internet. Mobile phones using WAP do not have as much screen real estate (the Nokia 9110 is shown in figure 3.9) so such users will probably prefer short and crisp versions of the documents when they use a mobile phone. Nokia is working with 3com to produce a hybrid between a Palm Pilot and a mobile phone.

This tendency might as a side-effect be beneficial for the currently rather overlooked blind computer users (see the Blind Web-ring³⁶ for more information) . They must normally use text-to-braille software or a “screen reader” in order to use computers, neither of which work well with web pages loaded with graphics, but when web authors must write for many kinds of media instead of just a particular version of a given browser, this will also be beneficial to these users since it will be easy to write a renderer for their preferred format. Of course, it is possible to write ordinary HTML in a way that is usable by these users, but that is rare these days.

A reference to the Jakob Nielsen site on web-accessability

³⁶ ...the Blind Web-ring... - <http://www.webring.org/cgi-bin/webring?ring=blind&list>

3.7 SGML/XML: Generic formats for storing data and meta-data

The problem of representing content in a generic way is not new at all, and was solved for the printing industry in the 1970'ies by the design of the Standard Generalized Markup Language (SGML) which is widely used. SGML was deemed too complicated for web browsers, and a trimmed version was named XML and standardized in 1998. Chapter 7 on page 53 talks a lot more about this.

For now it is sufficient to say that it is possible to store both data and meta-data in a convenient form in a SGML or XML file, but that these files must be rendered into the formats expected by the users, like HTML, PDF or what else tomorrow may bring of new, exiting formats. Since SGML describes the *content* and not the layout, it is “just” a matter of creating a renderer for any new format and add it to the collection.

This has been recognized by almost all of the “Documentation Projects” (The Linux Documentation Project³⁷, The FreeBSD Documentation Project³⁸ and others) accompanying the various free operating systems available on the Internet.

Chapter A.5 on page 99 talks about rendering documents to a given format on demand (for example in a web browser).

Should I talk a little bit about O'Reilly?

³⁷... The Linux Documentation Project... – <http://www.linuxdoc.org/>

³⁸... The FreeBSD Documentation Project... – <http://www.freebsd.org/docproj/docproj.html>

3.8 The user should use familiar tools to publish documents

SGML-editing is hard and tedious - this is one thing that the folks in the `comp.text.sgml`^a and `comp.text.xml`^b newsgroups agree upon. Tools are essential for authors.

Very few tools exist, and the good ones are quite expensive. The general consensus is that the best OpenSource tool is the Emacs editor with the PSGML package (see A.5 on page 99), and that this tool is primarily suited for programmers and other people who are well acquainted with SGML and XML.

PRINT DIRECTLY TO A PDF FILE - WEB PUBLISHING

Computers should *help* you do your work, and users are generally most productive with the tools they know. Why should an author use SGML with an editor she loathes, if almost the same result can be achieved by using a word processor she is familiar with?

In order to automate the conversion from e.g. Word to the equivalent SGML file, it is important that the author uses only those constructions that the conversion program understands. This is most likely given as a set of macros that must be used, as well as a verification program which the author can use at any time to check whether her document is conformant.

Publishing a document to the web should be as easy as printing a document. The basic principles are the same - you can do with a “Publish”-button and a dialogue box where the essential information is provided. It is not so today - discuss reasons . Use “print to fax” software as horrible example of this idea gone wrong.

When the webserver is dumb, you cannot do much. If there is a database underneath, much more is possible. Discuss the idea of having several ways of entering documents in the database (upload via form, send as email (fax software can do this too), print to virtual printer, scan, fax, voice) and letting the software do the conversions.

Use example with print PostScript to file and upload via ftp to printer (LexMark).

Writing XML directly is much harder to do than writing HTML (stricter syntax, more options, plainly just more to type), and should be aided by a good tool. Alternatively, the user should use well-known tools (Word) and mark up according to very strict rules, which is then automatically converted to the XML document. This does not give as rich documents, but allows users to publish existing documents with very little trouble.

^a...`comp.text.sgml`... - `news:comp.text.sgml`

^b...`comp.text.xml`... - `news:comp.text.xml`

Chapter 4

The need for well-defined standards

The need for open, well-defined standards. HTML, PDF, PNG (versus GIF) RFC-things (SMTP, POP3, IMAP, MIME) good examples. Word counterexample (look for Bill Gates saying people can buy an upgrade to read Word97 files in Word95). Adobe has good documentation on their PostScript and PDF document formats (doing theirs to have it well supported). HTML is widely documented in RFC.

4.1 Why is Open Source essential?

Independence from hardware and software vendors, allowing free choice of hardware platform and operating system, as well as the possibility to fix bugs and use other products. (Examples: XT versus Xerxes, Java Interpreter for Linux, servlet technologies, cheap - use money for hardware instead of licenses).

Chapter 5

Sample websites

slashdot.org (examine code - ask developer), Politiken, DynaWeb (SGI - oooold check on fyn, when was it first developed?), valueclick (banners - ask Ask for full story. What are they running), php3 documentation online. photo.net (reread). Sun's docs.sun.com (DocBook SGML)
Talk about standard search engines on web pages.

5.1 slashdot.org - high volume information site for nerds

Describe setup. Describe flow, and the slashdot effect. What hardware/software.

SlashDot¹ was started in September 1997 with the slogan “News for nerds - Stuff that matters”, and do so by providing dynamically generated web pages with a lot of stories categorized with an icon, where a short summary is shown along with links to the full story, several references to the original sources, the authors email address, plus an overview of the number of comments in the discussion forum. See figure A.5 on page 99 for the top of the start page at March 14 2000.

Slashdot Stats

```
date: 2:37am
uptime: 68 days, 5:56, 2 users,
load average: 0.18, 0.19, 0.18
processes: 65
yesterday: 209994
today: 1
ever: 278137148
```

Slashdot provide the “magnet content” Greenspun talks about. Stories are carefully selected, and the user may even select just a few categories in her personal preferences.

¹...SlashDot... - <http://slashdot.org>

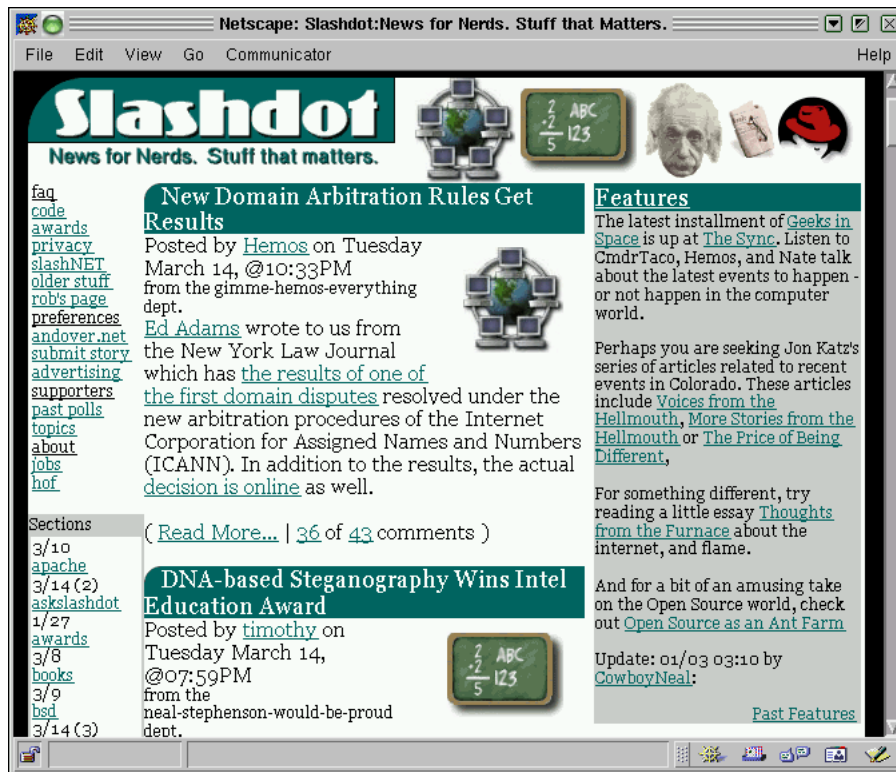


Figure 5.1: The Slashdot start page

All pages are dynamically generated. The discussion forums for any story is active for a reasonable time, where anybody can comment on any other comment, and then “frozen” when the story is archived. If a user at a later date wishes to see the story with the discussion, it is retrieved from the backup database.

The “Ask Slashdot” feature is a direct consequence of the fact that Slashdot is targeting power computer users. By allowing carefully selected questions to a highly visible spot, the original inquirer is usually able to get a very qualified answer. The question along with all the answers are stored in the same way as a story.

5.2 Valueclick.com

Banner provider. Check their online specs. Interview Ask for details, about how the NT-cluster was replaced with a small Linuxbox with MySQL and a squid in http-acceleration mode. Explain the problems with having fat processes serving slow modems.

5.3 Amazon - an Internet bookstore

The Amazon Internet bookstore² pioneered the virtual Internet bookstore, where the potential buyers use their browser to see the virtual inventory of literally millions of books, and order their selections. These are processed by Amazon and forwarded to subcontractors who deliver the purchased goods by mail to the buyers. The concept has since been copied by several others.



Figure 5.2: The frontpage of Amazon.co.uk. An ISBN-number has been entered in the very prominent search box

The Amazon frontpage is shown in figure 5.2. Note the high amount of potentially interesting links, and the prominently placed search box. Due to the limited domain of books, the search engine can make assumptions about what the user is requesting information about. In this case a ten-digit number was entered which happened to be an ISBN-number, which in most other search engines should have been entered as “ISBN-1-55860-5347”, and the page for the corresponding book was returned (see figure 5.3 on the next page).

This page shows the essential information like author, price, and the number of pages, but also an image of the cover, the rank on the Amazon sales list, and a direct link to purchase the book (Amazon remembers your credit card information, allowing for their patented “One-click-purchase” technology). What makes the Amazon web page truly better than a printed catalogue is that the users are encouraged to submit feedback, and the information Amazon collects from online sales to provide information to potential customers about other books that might interest them, because they interested other customers which bought the book they are looking at now.

²...Amazon Internet bookstore... – <http://www.amazon.co.uk>



Figure 5.3: The result of the search in figure 5.2 on the preceding page

These features are distinct for a page for a purchasable item at Amazon:

- **Average customer rating** - this is the essence of all the reviewers opinions, on a scale from 0 to 5.
- **Reviews** - readers take the time to create comprehensive reviews of the books, all of which Amazon then places prominently on the page, with due credit to the individual author. Recently a “Was this review helpful to you?” button was added to each review, and the result of each poll is shown along with the review, allowing a customer to rapidly determine how to rate the review.
- **Similar books** - if a previous buyer bought other books along with this one, they might be interesting to this person too. Links are presented to those books, as well as to the general categories this book was placed in by the Amazon librarian.
- **Expected content** - Amazon can be expected to have *any* English book in their database. Users expect to be able to find a given book in the Amazon database, and usually do. Personally I have never looked in vain at Amazon for English books.

Amazon explicitly invites those who read the page to review the book if they have read it, with special treatment for the author and the publisher. *As Greenspun observes, then this is the most efficient material*

At the time of writing, this book has 9 reviews at amazon.co.uk, but 198 at the mother site at amazon.com. The corresponding HTML-pages are 22 kb and 550 kb respectively, which gives an impressive *Greenspun factor* (setting the UK-version to be 100% *??-provider*) of 25.

Amazon uses very aggressive marketing, which has caused it to be a well-known brand name in just a few years. By giving money to those who link to their pages, as well as pay search engines to have links to the Amazon site show up at the top in just about every search, they have managed to be just about the only Internet bookstore known to most people. In section [vref\[??\]](#) I discuss why this approach ruins search engines for the ordinary users.

Conclusion: The Amazon family of sites are very good, and do what they can to improve with the help of visitors. Amazon have been able to keep their leader status in the virtual bookshop niche, by living up to the customers expectations. Hopefully, they will make money one day.

5.4 Bang and Olufsen

The Bang and Olufsen³ web-site was previously very unusual in the way that they offered navigation. Their front page contained very few links, and some product images which could be clicked and lead to a large set of pages (see figure 5.4 on the next page, with three random clickable images at the top. These were chosen randomly from the large set of such pages (600 or more) and put at the top. This has been highly efficient in keeping people wandering around and reading their pages, giving these people an unusual experience.

The underlying database is both used for storing the many pages, but also for tracking the users as they wander around the site.

³...Bang and Olufsen... – <http://www.bang-olufsen.com/>

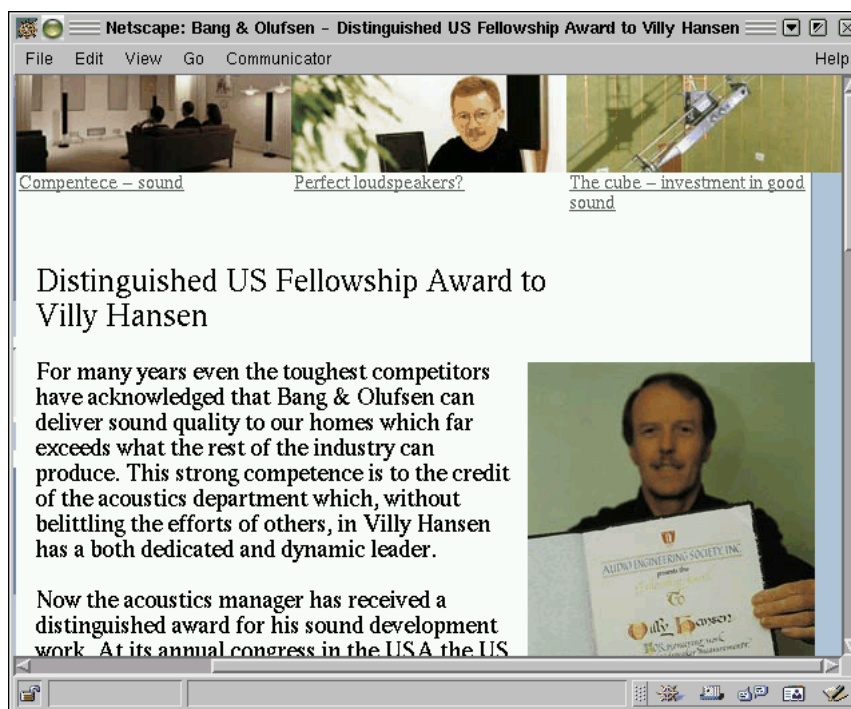


Figure 5.4: The Bang and Olufsen pages – notice the three randomly selected images at the top

5.5 Deja - where did Usenet go?

5.6 www.krak.dk

Making maps on demand was one of the earliest applications for the web. The *web map application was started* was discussed in section 3.1 on page 10, and has since been followed by *yahoo map, whoelse.*

In Denmark, the well renowned Danish map manufacturer Krak A/S⁴ opened their *web map in 199?* in direct competition with *909, 118 and who else?* since they provided a *person/number lookup facility*. Their advantage was up-to-date data *from the Telephone companies* combined with an ability to generate a map for a given address, whether it was entered directly or being a byproduct of a search for something else.

The website has been steadily improved with new facilities and cross-references, as well as more precise data (In *??* 1999 it could not locate “Herlev Hovedgade 205” on the map, so the whole of the road was inked on the map. This is corrected today). The start page is shown in figure 5.5 on the next page, and is separated in two parts; namely white pages (phone number based), and pink pages (company based).

The fields in the white page section are Name, Road, House number, Zip code, City and Phonenum-ber. The fields in the pink page section are Company and Phonenum-ber. The form has been filled

⁴...Krak A/S... – <http://www.krak.dk>

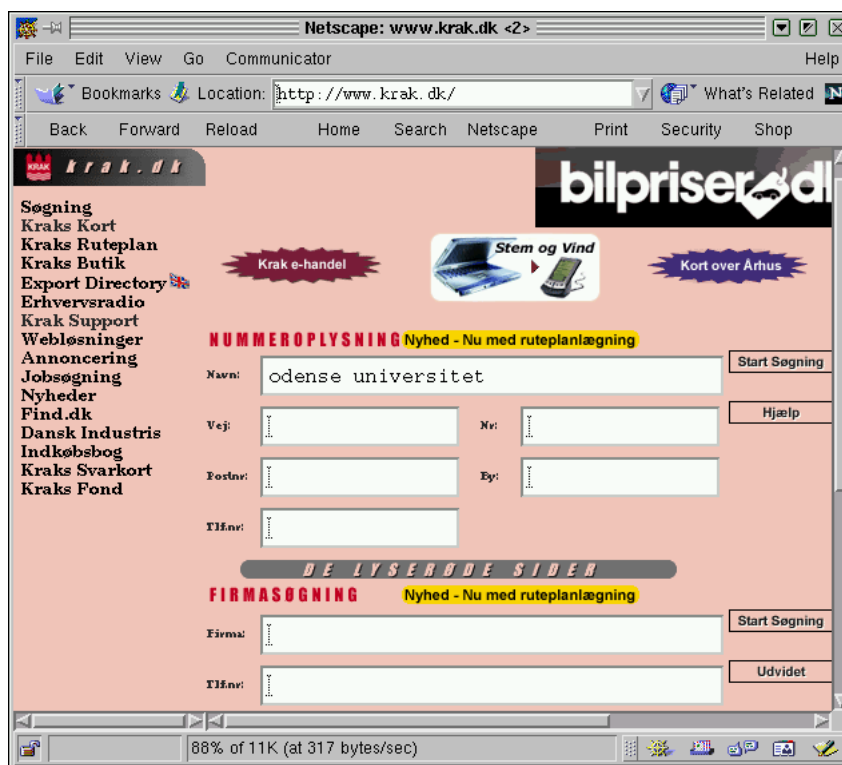


Figure 5.5: Welcome page (and search form) for www.krak.dk

out with a search request for “Odense Universitet”, and the results are shown in figure 5.6 on the following page.

Each line containing an answer may have icons referencing to facilities regarding the location of the answer. These answers have pointers to Rejseplanen, the Route Planner, and a map reference. Additionally pointers to a web page, and an email address could have been provided by the users. Clicking the map icon of the first line actually referring to Odense University brings us to figure 5.7 on the next page.

The blue dot usually indicates the location of the address but in this case the map is not 100% correct (Moseskovvejen does not go all the way south-east to the parking lot, and the University is located 500 meters further to the south). When a map is display, the navigation area below the image allows for zooming and moving the contents of the shown area. By selecting “Zoom out” and “x8” and clicking on the blue dot a new map is shown (see figure 5.8 on page 37), where it is evident that the database show a great deal of detail. Roads, residential areas, streams, train stations, the motor way and green areas are shown. These maps provide a level of information corresponding to Krak’s printed maps.

Now go back to the list of results from the search for “Odense University”. The Car-icon gives a route to the listed location, where you must fill in the source yourself. Figure 5.9 on page 37 show the form filled in with my own address and the address of the university, and “Route on map” (Rute på kort) gives a visual route between the two locations, along with an estimate of the distance and time it will take.

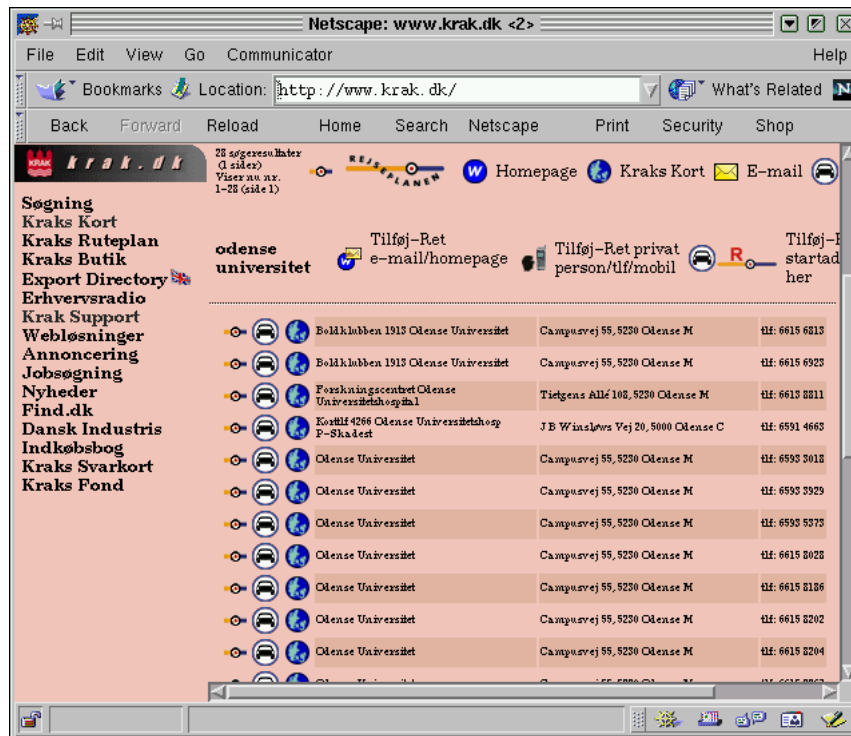


Figure 5.6: Result of searching for “Odense University”

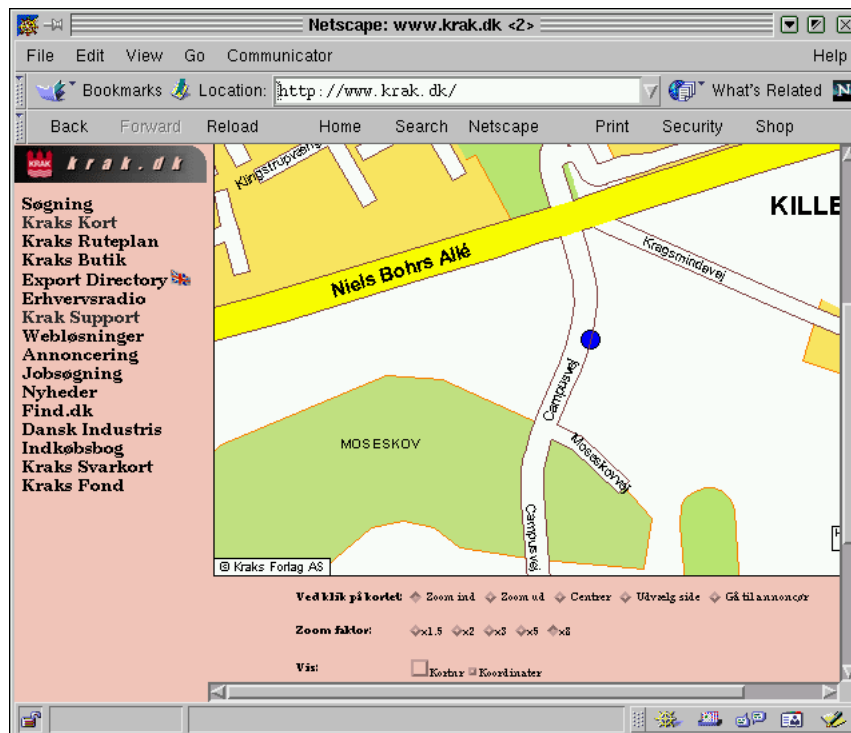
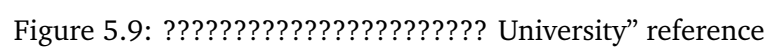
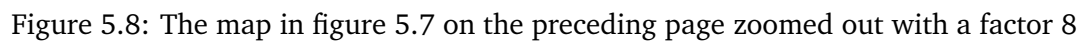


Figure 5.7: Following the map icon for the first “Odense University” reference



This map is shown in figure 5.10, and is perfectly adequate for a person driving a car. For cyclists it is often a good idea to look for short-cuts, if you are well known in the area.

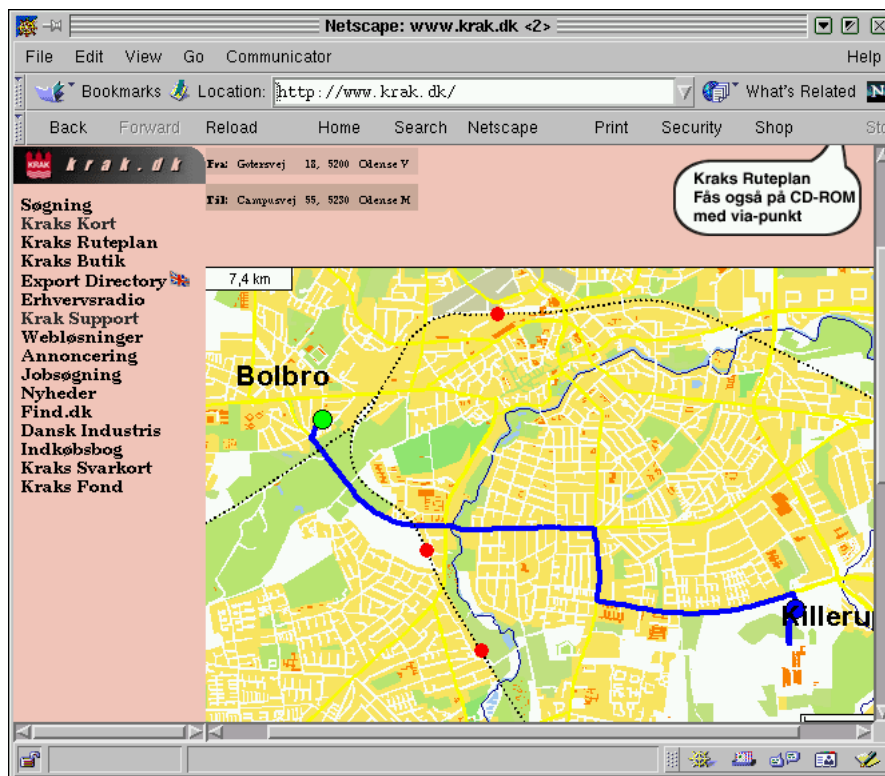


Figure 5.10: "University" reference

If you need step-by-step directions, Krak can provide that too. Select the "???" button and print out the directions. Careful studies of these directions can give an idea of the network grid that the route planner use for finding the shortest road. *what is 5.11 on the facing page?*

5.7 The Journey Planner - www.rejseplanen.dk

DSB (Danish Rail) have a journey planner which is based on train stations, and major bus stops. Figure 5.12 on the next page shows the initial form where the two end points for the journey as well as the arrival or departure time is indicated. Please note that the two surrounding frames (left side and top - the scroll bare on the right indicates the actual area available) leave only *70% of the area* to the application itself.

The search returns a number of potential journeys shown in figure 5.13 on page 40, with departure and arrival times, total expected travel time, and the number of changes necessary during the journey. The departure corresponding the best for the indicated period is highlighted with the blue bar. By accepting the default selection in the blue bar, and scrolling down to press a button, figure 5.14 on page 41 is shown. This is a direct connection, and it is possible to reserve a seat directly from this screen.

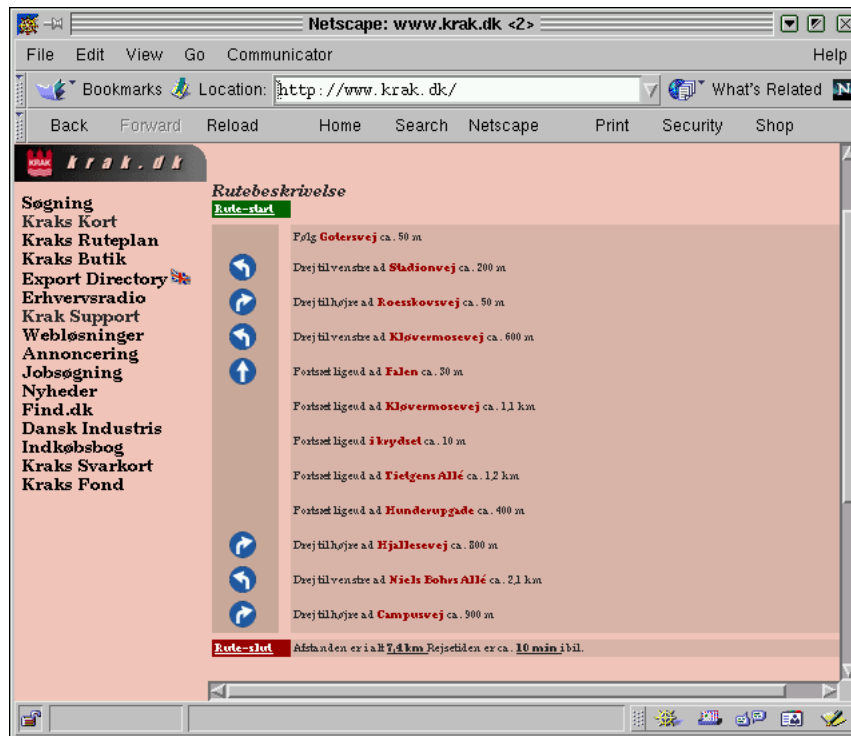


Figure 5.11: "University" reference

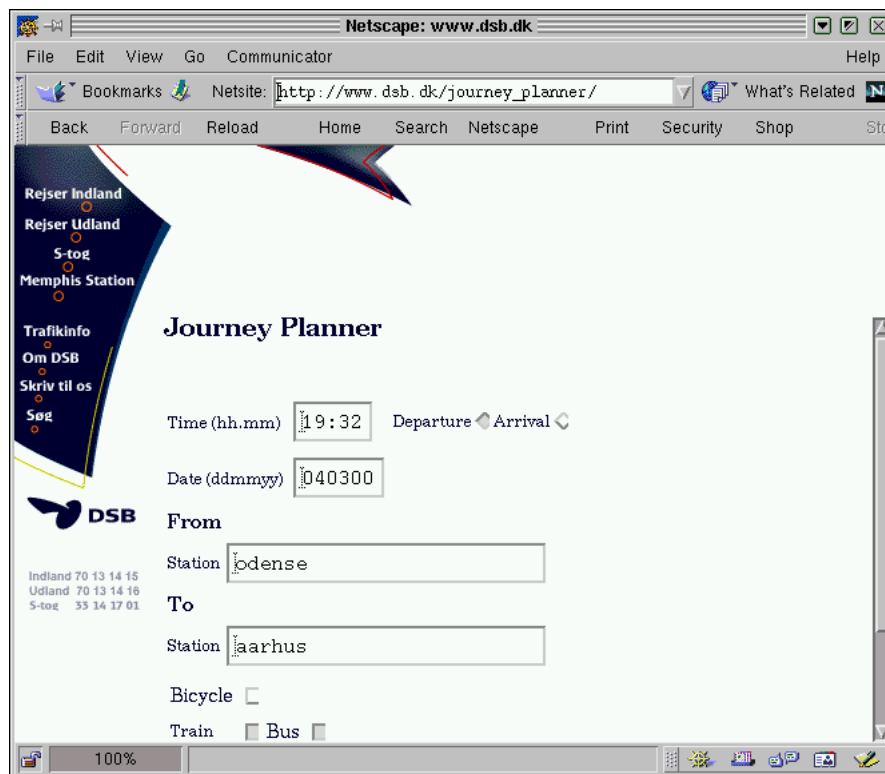


Figure 5.12: Rejseplanen1

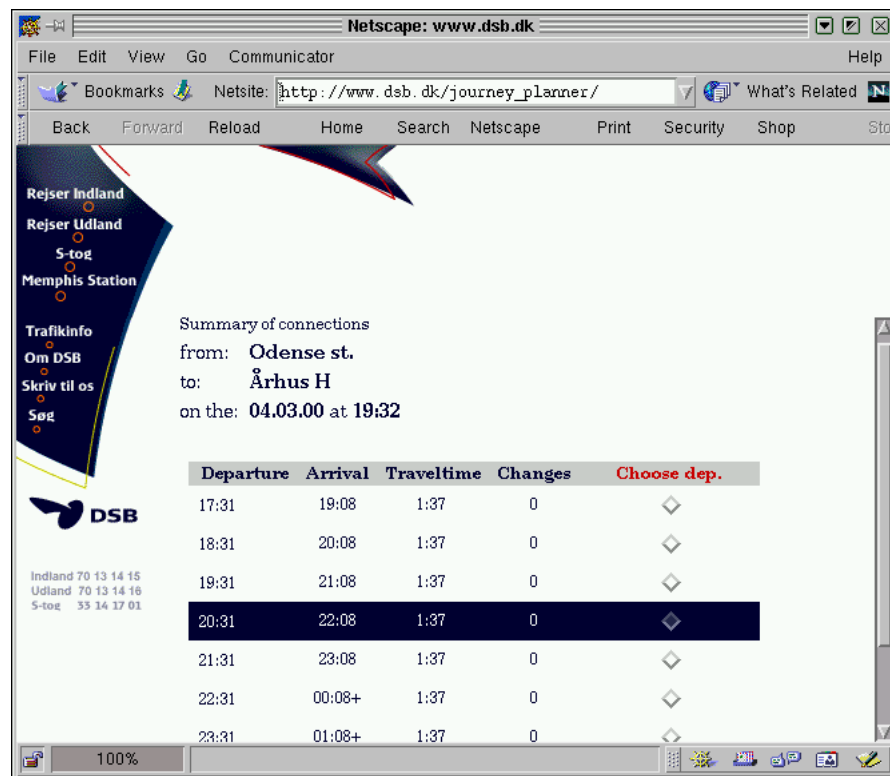


Figure 5.13: Rejseplanen3

A more complicated request is shown in figure 5.15 on page 42 where a connection from Malmparken (suburb to Copenhagen) to Esbjerg (Western Jutland) is requested. A route involving three trains and a bus is suggested, which is very reasonable and probably the fastest too.

Rejseplanen is helped by the fact that the number of nodes in the grid can be small compared to `www.krak.dk`, but the presentation is not good enough. It is generally too hard to get to the information if you need a slightly alternate view from what the programmers expected.

I used to live next to Odense Sygehus, which is a stop on the Odense-Svendborg railroad, but where not all trains stop. The question regarding whether it would be faster to walk to the Odense Sygehus station or take my bicycle to Odense Station was very hard to answer by Rejseplanen, since it does not just simply allow access to the underlying database. I ended up calling the station and asking them to look in their paper copy.

5.8 Freshmeat.net

Freshmeat.net is a service for locating software. *a lot of products. users can announce software, and point peo*

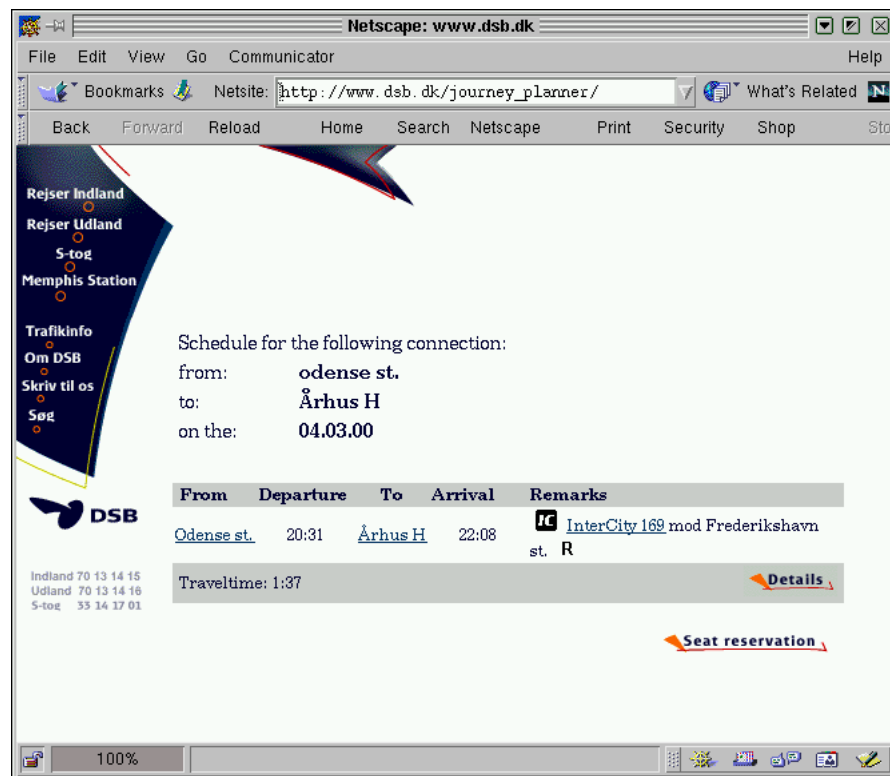


Figure 5.14: Rejseplanen4

5.9 Altavista - the search engine

AltaVista was the first web site to proudly present a search interface to *all* the web pages on the Internet⁵, where they extracted words from the individual web pages to create their search indexes. They gave extra weight to the descriptions that the web page authors entered themselves in the META-tags (see A.5 on page 99), in order to present the “hits” sorted with the best hit at the top.

1. Words in META-tags
2. Words in titles and headlines
3. Word in the body of the text

Check if there still is companies which boost your presence on web search engines

Businesses - notably in porno - discovered this and polluted their web pages with keywords to boost their

Look for statistics on search engines and “sex” - talk about that these businesses have a great interest in a

look for email about somebody who tried google and was pleased after being driven away from search en

⁵Later it was revealed that AltaVista only indexed up to 50 pages for a given site, meaning that big sites simply weren't indexed. That gave the other search engines like Lycos and **Inktomy** a parameter to compete on

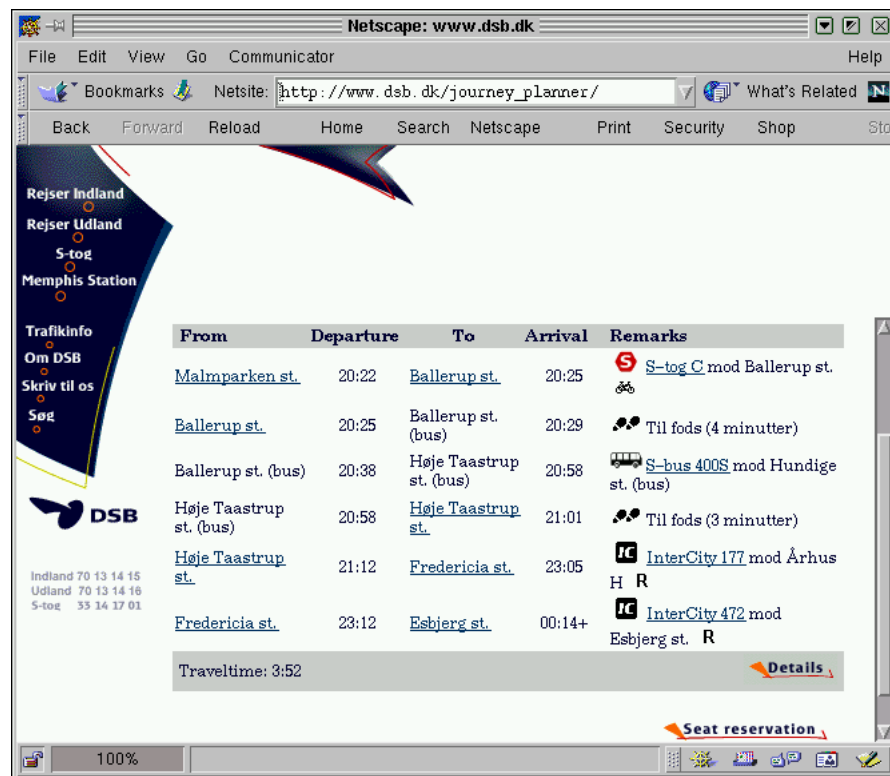


Figure 5.15: Rejseplanen5

ftpsearch.lycos.com⁶ is a unique service provided as a byproduct of the main search engine at www.lycos.dk⁷ where you may search for *images* on the web. The search engine does not look inside the images it finds while traversing the net, but looks at the filename and the context in the web page for hits. It works reasonably well - most of the cactii images for my project was found this way. Check on details how they do it.

Search engines. How do they do it? www.909.dk, www.krak.dk, ter-raserver.microsoft.com. Encyclopedia Britannica. ASP+Access, tinderbox & bugzilla, javasoft (developers area), LXR+Bonsai. "Alle danskere paa nettet" - Dansk Journalistforbund, Pressemeddelelse. Greenspun selv?

5.10 The Collection of Computer Science Bibliographies

<http://liinwww.ira.uka.de/bibliography/index.html>⁸

screens hot

⁶...ftpsearch.lycos.com... - <http://ftpsearch.lycos.com>

⁷...www.lycos.dk... - <http://www.lycos.dk>

⁸...<http://liinwww.ira.uka.de/bibliography/index.html>... - <http://liinwww.ira.uka.de/bibliography/index.html>

This is a great web-site for scientists because in addition to the basic search facilities it allows you to download the BibTeX entry for the books and articles you find.

5.11 Jydske Bank

The only java-solution for home banking. Any comments?

5.12 Cactus – document capture and conversion

Two out of four methods are implemented (email, printer). (Increase blob limit in DBI::MySQL). Sample PDF viewer, perhaps simple XLS viewer.
Demonstrate that Cocoon can provide the navigational framework.

Chapter 6

Overview of technology available for Linux February 2000

Describe what products there are currently available for these categories. Evaluate as much as possible. Conclude that currently this is an area in great development.

6.1 Webservers

Apache, misc Servlet Javasers (jetty, thttpd, java webserver, etc). Look for comparisons. Roxen (graphs). Discuss whether several webservers should be run on the same machine to provide better services.
Look at the netcraft search - differentiate between Internet server and intranet server, and internal product stuff.

6.2 Database engines

The number of database vendors which support Linux (which as for now is *the* Open Source operating system which people know about) is steadily increasing. This webpage try to keep track of them all¹, but I only list the major commercial vendors and Open Source projects, where the drivers should be of production quality. The following overview was created March 2000.

PostgresSQL (discuss features: transactions, unique numbers, in-kernel code, speed, availability, multi-host capabilities?). Any flat file RDBS?
Greenspuns evaluation of embedded webservers in databases

¹... This webpage try to keep track of them all... - <http://linas.org/linux/db.html>

6.2.1 Cloudscape - Informix

The Cloudscape Database² is written in Java, with a JDBC and a HTML interface, has many core SQL and Java extensions implemented, and allow distribution amongst several Java-machines (using RMI). Informix bought Cloudscape Inc. in October 1999 to get this database system, so it is still very new. They offer a 60 day trial version from their website³. Pricing is approximately \$900.

6.2.2 DB2 - IBM

DB2 for Linux⁴ is available - the <http://www6.software.ibm.com/dl/db2pde/db2pde-p>⁵ is available for free for non-commercial purposes.

DB2 XML Extender⁶ is an interesting addition to DB2 which “let you store XML documents in DB2 databases and new functions that assist you in working with these structured documents.”. It is available for AIX, Solaris and Windows NT.

6.2.3 Informix - Informix

<http://www.informix.com/datablades/dbmodule/informix1.htm>

6.2.4 Ingres - Ingres

The Ingress II database is in a beta stage for Linux⁷, and can be downloaded for free. It uses Perl, gcc and Apache to provide webserver facilities.

6.2.5 mSQL - ?

Yes? Looook for stuff - rememebr to get the mSQL/MySQL book in the databse

6.2.6 MySQL - TCX

MySQL is an Open Source database⁸ which is very popular amongst Perl programmers due to the good performance⁹ combined with the free availability of source code and high quality drivers.

MySQL only costs money if you run it on a Microsoft operating system, or commercially on a server.

²...Cloudscape Database... - <http://www.informix.com/informix/press/1999/dec99/cloudscape.htm>

³...a 60 day trial version from their website... - <http://www.cloudscape.com/Evaluations/index.html>

⁴...DB2 for Linux... - <http://www-4.ibm.com/software/data/db2/linux/>

⁵...<http://www6.software.ibm.com/dl/db2pde/db2pde-p>... - DB2 Personal Developer's Edition V6.1

⁶...DB2 XML Extender... - <http://www-4.ibm.com/software/data/db2/extenders/xmltext/>

⁷...Ingress II database is in a beta stage for Linux... - http://www.cai.com/products/betas/ingres_linux/ingres_linux.htm

⁸...MySQL is an Open Source database... - <http://www.tcx.se>

⁹...the good performance... - <http://www.tcx.se/benchmark.html>

6.2.7 Oracle

The complete Oracle 8i product range is available for Linux, with a low-end version available for download. The WebDB tool functions as a webserver-interface to any Oracle database.

When Oracle visited FLUG¹⁰ we were promised that we could have all the Oracle software for Linux we needed (without official licenses), and that Oracle Denmark would provide the necessary support. It is interesting that Oracle Denmark unofficially seed the Linux communities with pirate copies of their software.

6.2.8 SyBase - Sybase

Sybase has Sybase Adaptive Server Enterprise and SQL Anywhere Studio¹¹ available for Linux for free as long as they are used for development. For production installations a license must be bought. An earlier version is unsupported but may be used without restrictions for production environments.

Their web integration tool does not appear to have been ported to Linux yet. The Linux database page¹² lists several third-party tools which provide web functionality.

Even though Microsoft does not provide Linux drivers for the Microsoft SQL-server, the freely available SyBase drivers for Linux appear to be usable.

6.2.9 Other approaches to database storage

Landbrugets rådgivningscenter¹³ uses a Lightweight Directory Access Protocol¹⁴ server to provide easy access to commonly used web elements.

6.3 Browsers

MSIE5 (XML support, Solaris), Mozilla (alpha), Netscape (too old), HotJava 3.0 (no XML, can be expanded to do XML conversion internally?), Opera (?), Amaya (other W3C browsers?), AWTviewer in fop (can it do xml directly?),

6.3.1 Netscape Communicator

The Netscape Communicator browser is available and fully supported for Linux. The version available at March 2000 was 4.72 in English only.

¹⁰ ... Oracle visited FLUG... – http://www.flug.dk/pages/foredrag/tidligere_foredrag/tidligere_foredrag.html

¹¹ ... has Sybase Adaptive Server Enterprise and SQL Anywhere Studio... – <http://www.sybase.com/products/linux/>

¹² ... Linux database page... – <http://linas.org/linux/db.html>

¹³ ... Landbrugets rådgivningscenter... – <http://www.lr.dk>

¹⁴ ... Lightweight Directory Access Protocol... – <http://www3.innosoft.com/ldapworld/ldapfaq.html>

This browser have a reputation of being rather picky of its environment. Especially the Java-subsystem is prone to crashing the server. If Java is important, consider HotJava.

6.3.2 Mozilla - Netscape 6.0

The Mozilla browser¹⁵ has been under development since 1998, where Netscape – inspired by the *the Cathedral and the Bazaar* paper – decided to switch development strategy for its incomplete source for Netscape Communicator 5.0 from an in-house system to a full-scale OpenSource model housed at <http://www.mozilla.org>.

In the past two years, the OpenSource developers have basically rewritten the Communicator from scratch into a truly impressive product (this is based on the M14 milestone release). A separate beta-branch of the source code designated for the Netscape release was opened March 2000. A tentative release is midsummer 2000.

Mozilla is not yet suited for mainstream use.

6.3.3 Opera

This Norwegian browser¹⁶ has been highly successful on Windows for providing a small and quick browser which goes well with users who think that Netscape and Internet Explorer are too big and slow for their machines. Their Linux port is still in Alpha-testing, and unsuited for mainstream use.

6.3.4 HotJava 3.0

HotJava is a Netscape 3.0 compatible browser written in Java.

It is a bit slow, especially in screen updates, but has its force in its ability to execute Java applets (which is where Netscape is slow, and Microsoft does not comply with the standards).

HotJava requires a Sun JDK-based Java implementation (clones do not provide all the functions they need), and I have tested version 1.1.2, 1.1.4, and 3.0 under Linux, OS/2, Windows, Irix and Solaris, where it worked satisfactory.

There is currently no information regarding the status of parsing XML internally. Full access to the source is probably necessary before any third-party implementations come. In a perfect world, Sun would release the source, but that is probably very unlikely.

Rumours has it that development on this browser has ceased. I tried to extend version 1.1.4 to show TIFF files, but had to give up due to the miserable documentation for HotJava combined with lack of the source.

Other rumours has it that Sun will release the source for HotJava under their Community License (which basically makes the developers unpaid employees of Sun) like it has been done previously. On March 13, Sun released their first product under a true OpenSource license (the Forte product -

¹⁵... The Mozilla browser... – <http://www.mozilla.org>

¹⁶... Norwegian browser... – <http://www.opera.com>

previously Netbeans), which could indicate that they had evaluated their experience with previous releases, and found that nothing short of full OpenSource will do.

Remark: Today there are so many OpenSource projects that the truly talented coders can choose freely what they want to develop on, and they choose projects where they can get return in recognition from their equals, and where they don't feel that they slave for a company.

Sun has a "vote for the bug you want fixed"¹⁷ system for Java, and the highest rated bug was "Port Java to Linux" for a very long time, and with four times the number of votes for number 2. The Blackdown Team¹⁸ ported Java 1.1 to Linux, and collaborated with Sun on porting Java 1.2 to Linux. It took very long time, and when they finally had it officially published on the Sun Java webpages, Sun did not give credit to the Blackdown Team - not on the web pages, not in the code [rumours have it they even removed comments crediting Blackdown]. The Blackdown Team was highly offended, said publicly so, and stopped working on the port!

My personal guess is that this incident taught Sun that they must be very observant to the OpenSource culture they want to attract. It must be nurtured and credited in order to create the symbiosis that the Mozilla project has demonstrated to be possible.

In the same spirit SourceForge¹⁹ are providing free hosting for OpenSource projects, and they have reached the point where people ask "why don't we put this on SourceForge to get CVS".

6.3.5 Microsoft Internet Explorer 5.0

Microsoft has created a very nice browser which works well on Windows platforms, and which supports a recent version of XML with the XSL-stylesheets ([reference](#)). Since the XSLT standard was released January 2000, a release of MSIE50 can be expected "rather soon" which implements the XSLT standard, allowing the browser to render XML files on the client side. For now, server side processing is still the only option.

Microsoft Internet Explorer 5.0 is *not* available for Linux. It is, however, available for Solaris and HP-UX so the codebase is ported to Unix, and from there it is usually comparatively easy to get software to run under Linux. The Solaris version was incompatible with the window manager I was running under Irix, so I have not tested it much, but it was very, very similar to the Windows version.

My personal guess is that Microsoft holds back Linux versions of the Internet Explorer with Outlook Express until they see a fortunate opportunity to release it. A certain lawsuit springs to mind²⁰ - I do not expect Microsoft want to be accused of monopolizing the Linux market too.

To summarize: Microsoft Internet Explorer is not available for Linux at the time of writing.

Try Solaris version of MSIE under KDE

¹⁷... "vote for the bug you want fixed"... - <http://developer.java.sun.com/developer/bugParade/top25bugs.html>

¹⁸... Blackdown Team... - <http://www.blackdown.org>

¹⁹... SourceForge... - <http://www.sourceforge.net>

²⁰... A certain lawsuit springs to mind... - <http://www.idg.net/microsoft>

6.3.6 Amaya

Amaya is the W3C testbed for HTML-development. Is a nice browser combined with a reasonable editor. Version 2.4 do XHTML and HTML-4.0 with CSS but not XML yet.

6.3.7 StarOffice

What can we do here? - try it out.

StarOffice include a Netscape 3.0 compatible browser, with Java applet support. It is an integrated part of the StarOffice program, which besides Linux is available for Windows, OS/2 and Solaris.

As Sun recently bought StarDivision who had written StarOffice, the software currently being re-targeted to Sun interests. It is not evident yet what that might be, but it is not unreasonable that the thin Java client (connecting to a StarOffice server on a Solaris machine) is what Sun bought it for.

Currently there is no sign for integrating XML support in StarOffice.

6.3.8 W3 - browsing inside Emacs

The W3-browser for Emacs²¹ (needs XEmacs to show images) is a reasonably good browser, which has its force in its tight integration with Emacs, and is an interesting alternative for those who live in Emacs anyway. W3 is written in Lisp. The next major version is expected to be rewritten around a generic XML parsing engine.

6.3.9 Lynx - browsing in textmode

Lynx is one of the oldest browsers around still in widespread use, and does text-mode only. It is being steadily improved and hacked on, both because it is a nice browser for those who may not have a graphical login available, but also because it has a lot of options to help the programmer debug web servers and scripts. Has only recently been superseded in usefulness by the various tools in the Perl community.

²¹ ... W3-browser for Emacs... - <ftp://ftp.cs.indiana.edu/pub/elisp/w3/>

6.3.10 Perltools?

6.3.11 OTHER BROWSERS?

6.4 XML utilities

XT, Xerces, Cocoon, SQL2xml, validators (error recovering?), converters to XML (pod2docbook, tex4h, Perl with appropriate modules?, wordview modified)?
Show the XML- \rightarrow SQL- \rightarrow XML- \rightarrow HTML on the fly conversion possible with Cocoon Without auxillary code. All done in W3C style.
WordPerfect 9 should be able to edit SGML directly.

6.5 PDF utilities

PDF generated by pdflatex is good. PDF generated by fop isn't. createpdf.adobe.com
pdfzone. acroread (ok, heavy), xpdf (font problems with bitmaps, ugly but fast, gives best result with usepackage times, no anti-alias), gv (pretty good). Distiller (check what it can do). Use pdftotext to get the number of pages in the pdf file.
Conclude that it is possible to have a 100% Java solution for on-the-fly publishing XML to HTML and PDF through a webserver. Discuss why developers do it in Java, as opposed to any other language.
Very few utilities support creating PDF. Discuss how to create PDF.

Chapter 7

SGML, XML and DTD's

A major ingredient in developing large, complex systems with several layers of data and meta-data, is the ability to *abstract* between different levels of the system.

- Operating systems understand the “file” abstraction, so the programmer does not have to think about the underlying hardware when reading and writing information.
- Expensive printers provide the PostScript language which is an abstraction for “electronic paper”. Any Postscript file can be printed to any PostScript printer, without considering the specifics of the device.
- Programming languages, like C, provides the abstraction of a “virtual machine”. The programmer does not have to worry about how large an integer, and how the best way to implement a loop on this particular processor is. The compiler takes care of that, leaving the programmer to concentrate on problems at a higher abstraction level.

The same thing happened in the printing industry, where SGML (Standard General Markup Language) was developed in the late sixties and early seventies¹ to facilitate a transfer from “specific coding” to “generic coding” (Goldfarb mentions using “heading” instead of “format-17”). This allows authors to concentrate on *what* they write, instead of *how* it is presented. In Goldfarb's own words in 1971:

The principle of separating document description from application function makes it possible to describe the attributes common to all documents of the same type. ... [The] availability of such 'type descriptions' could add new function to the text processing system. Programs could supply markup for an incomplete document, or interactively prompt a user in the entry of a document by displaying the markup. A generalized markup language then, would permit full information about a document to be preserved, regardless of the way the document is used or represented.

Today, almost 30 years later, this is coming to the users of the world wide web, as the presenting technology moves into the web browsers - with XML being defined as a subset of SGML - and they will be able to fully present XML documents individually tailored to the users preferences. This is

¹... the late sixties and early seventies... – <http://www.sgmlsource.com/history/roots.htm>

a radical change from the current situation where HTML allows an information provider to be so specific in markup that it makes it unusable in some browsers.

When XML is common place in browsers, we will hopefully have much more content in webpages allowing for

7.1 “HTML is a SGML DTD” - the concepts

Table 7.1 shows how the most frequently used acronyms relate when mapped to the programming domain. See section 2 on page 5 for the full descriptions.

Documents	Programs
SGML	Family of programming languages (C/Pascal/Lisp...)
XML	Another family of programming languages. (XML is a subset of SGML)
DTD	Backus -Naur notation of a given language, defining how it looks
HTML	A certain programming language like C
XHTML	Another programming language like Java Well written programs can run both in C and Java
DSSSL	A report generator which takes a XML program as input and produce a report, another program or something else
XSL	A report generator which takes a XML program as input and produce a report, another program or something else

Table 7.1: The analogy between programming languages and SGML

Basically SGML (Standard General Markup Language) covers a large family of document markup languages, which all share that a given SGML-document *must* comply to a DTD (Document Type Definition). DSSSL is used to convert a SGML document into something else. DSSSLs are written in a Scheme like language (see [Dybvig(1996)] for a reference manual).

XML is a light-weight version of SGML designed to be embeddable in a browser. XML-documents *may* comply to a DTD, or be stand-alone (“DTD-less”). XSL is used to convert a XML document into another XML document (XSLT) or into the generic page description language XSL-Formatting Objects (XSLFO).

XHTML is a XML-version designed to be parsable by HTML-4.0 browsers.

7.2 The author should provide content, not do layout

SGML came originally to because the authors at IBM were not consistently marking their documents up, resulting in suboptimal documents. This was because they were doing physical markup instead

of logical markup.

Modern word processors have been constructed around the “What you see is what you get” paradigm, where a great deal of effort has been spent on allowing the user to see exactly on screen what would be printed on paper. Unfortunately the limited resolution and size of a modern computer monitor cannot show the whole page in high resolution, and the user must therefore work in small “regions” without being able to see the whole page, or perhaps even several pages A tufte reference here?. Microsoft Word provides a draft mode where just the text is shown without markup; WordPerfect provides a “Reveal Codes” mode where the user can edit the text along with the tags.

The most recent versions of Word and WordPerfect wp9 should have sgml builtin have had styles built in, where the user could specify “heading” and similar categories for selected text, but have been very passive in helping users to employ these categories, meaning that most of these documents have had almost no semantic markup, but only visual markup.

In combination with this the web publishing tools by Microsoft and others have blatantly ignored web standards, and used all kinds of font modifying tags instead of the HTML-tags conveying *meaning* about the text. These tags often implement the particular settings that the author had on her word processor, so that the author is effectively doing layout on the web.

This is not a new phenomenon, however. The best way to typeset mathematics have been the TeX system², (see section refsec:tex-and-latex) where the author basically is told to write what she wants to and let TeX do the rendering. Even so, it is often very tempting even for the most dedicated purist to want to change the layout to make a small visual change, simply because they *can*. Unfortunately very few have the expertise to do this right, so the changes usually make it worse. I have seen a book where the author had prepared camera-ready copy for his \TeX -document. The book was printed in A5, but he had handled this merely by increasing the margins of a normal page, and nothing else. The result was that the layout looked awkward - the line spacing seemed too large for the page. This would not have happened if a professional layouter had produced the book!

A layouter may also want to change the whole layout of a book, for example if a new edition is to be part of a series with another distinct style. The publisher may want to turn the book into webpages or a CD-ROM, or into a Braille version for blind persons. It might even be read aloud by a computer.

It should not bother the author during the writing process how the finished document will look, as long as she marks up the document appropriately, providing sufficient meta-data about the content.

As the possible annotations are decided by the DTD, it is important to choose a DTD carefully. A well annotated document can be used for many purposes, not all of which was envisioned when the document was written. Ten years ago no one had heard about the world wide web - within the next few years, it will probably be the main use of SGML in the world.

7.3 The creation of HTML

Look up creation of HTML on the net. Look for reasons why HTML was chosen to be a SGML DTD instead of

²...TeX system... – <http://www.tug.org>

The w3c have notes and DTD's describing the first versions of HTML³.

HTML: 1.0 (original, simple, <hr> tag was added due to popular request), 2.0, 3.0 never made it, 3.2 tables, 4.0 w3c straightens up things. Netscape added new tags *ad hoc* often. In the mean time HTML has been made to do things it was never meant to do originally, for presentation purposes (large imagemaps in tables without any textual information), and there is still a need for new features that designers want. *W3C saw that blindly adding new features to HTML would result in an even more bloated standard* with a lot of backward compatability to maintain - a modern browser must still be compatible with the bugs in Netscape versions 2 and 3. These were renowned for being very lenient towards users providing erroneous HTML, and have started a tradition of browsers doing their best to interpret what the users *might* have meant. This makes it hard to use HTML as a generic information data format, since parsers are complex due to the many exceptions.

The W3C took a bold step in saying: We need a new format - the eXtensible Markup Language. XML!

7.4 The creation of XML

XML was as HTML built on SGML, but with a much stricter syntax than HTML and a lot less facilities than in SGML. XML is a new standard - it was made a W3C recommendation in 1998, and the accompanying XSLT and XPath in November 1999.

What was the design goals

The resulting language for designing languages has these features:

- A very strict syntax which every XML-document must comply to (which is what a *well-formed* document does). All open tags must be closed explicately. Tag attribute values must be in quotes.
- Easy to parse and generate. The strict syntax makes the parser simple.
- The DTD is not mandatory. DTD-less documents does not have a DTD to conform to, and must only be well-formed. Unicode?

TRANSFORMATON

FORMATTING

7.5 XHTML - XML compliant HTML

[The XHTML W3C recommendation⁴.]

³...The w3c have notes and DTD's describing the first versions of HTML...– <http://www.w3.org/History/19921103-hypertext/hypertext/WWW/MarkUp/MarkUp.html>

⁴...The XHTML W3C recommendation...– <http://www.w3.org/TR/xhtml1/>

A big problem with XML in general is that it cannot be viewed by anything the majority of users have today. It is possible, however, to specify a form of XML called *XHTML* which bridges the gap between the two worlds of XML and HTML. Documents corresponding to a XHTML DTD are parsable by HTML-4.0 compliant browsers if a few, simple guidelines are followed:

- Elements (the HTML tags) must be closed correctly.
- Attribute values must be quoted
- Element names and attributes must be given in lower case
- Empty elements must either have an end tag, or the start tag must end with “/”>. I.e. a horizontal line is written as “<hr />”, where the space is important to make this acceptable to HTML-4.0 parsers.

The well-formedness of an XHTML document is achieved...

HTML: <hr noshade>

XML: <hr noshade="" />

Why is this smart??? Not only pure output but also input for another XML transformation

XHTML is an XML-variant designed to be viewable in HTML-4.0 compatible browsers (without XML support). XHTML-documents are well-formed XML documents, while at the same time being valid HTML-4.0, so that XHTML can be the target of a standard XML-transformation as well as a source for initial XML-processing.

This is very different from the usual situation, namely that HTML is the result from *formatting* a document, where the result is not XML any more.

Note: XHTML style sheet generated files does not trigger correct change to UTF-7 character set.

WHICH CONVERSION UTILITIES ARE THERE? TIDY? READ XHTML? WRITE XHTML? WHAT DO W3C SUG

7.6 Which dialect of SGML should be used?

The Document Type Definition (*DTD*) defines exactly the way a SGML/XML document can look, when conforming to the DTD.

- which tags are valid
- which tags can appear at any given part of a document
- which attributes are valid for a given tag
- which *entities* (macros) can appear in a document, both for expansion strings but also for Unicode characters unavailable to the document author
- **Others?**

There is a lot of different DTD's available to many different purposes, since basically everywhere a datafile is needed XML can be used with a suitable DTD. [ChemML](#), [SVG?](#). If a need arises it is just a matter of creating a suitable DTD to work with it.

Unfortunately, for each and every pair of DTD and output format an XSL-style sheet must be written, tested and maintained. Therefore it is a very good idea to use a commonly used, documented, and well-thoughtout DTD for the documents, and several other groups have already done such work.

The Cactus system uses the DocBook XML V3.1.7 DTD with the DocBook XSL [DocBook XSL](#)

7.7 How can a *ML document be viewed?

In order to be rendered into a view, a *style sheet* must be applied to the document. Style sheets convert the tags in the document to a given output format, and depends on (DTD, output format pairs).

GRAPH DEPICTING STYLE SHEETS BEING APPLIED FOR DIFFERENT OUTPUT FORMATS

Commercial products usually have a viewer which can apply the style sheet directly examples? Panorama? Framemaker+SGML, but are due to licensing costs and supported user base rarely an option for general for documents targeting end-users, and I have not tried any of them. Table 7.2 discusses the current *de facto* document formats on the Internet.

Format	
HTML	Browsers are available for any modern platform. HTML-4.0 compliant browsers are Netscape Navigator 4, Internet Explorer 4... More?
PDF	Adobe Acrobat Reader is available for most mainstream platforms URL . The Adobe Acrobat Viewer is available for any platform with Java URL . The Open Source project Ghostscript URL is available for these and other platforms.
Microsoft Word	The doc and rtf file formats are widely used, but needs a full word processor to format properly. Anything else but Microsoft Word gives inferior results, and Word only runs on Windows.
PS	
PDF	

Table 7.2: Frequently encountered document distribution formats

The Microsoft Internet Explorer 5.5 (not yet available at the time of writing) will be able to transform XML documents with XSL-style sheets internally complying with the W3C recommendation. The previous version of IE was available on Windows, Solaris and HPUX.

The Mozilla browser in its pending release as Netscape Navigator 5.0 will not be able to apply XSL-style sheets, but only CSS cascading style sheets. This [decision is rather old and may have changed](#).

For the immediate future web publishing in XML implies a need for server side conversion to other

formats. The available Open Source software for this is surprisingly small, but that is rapidly changing. Please note that XML style sheets still implies SGML style sheets, which is interesting because these are more mature due to longer use (see section~refsec:docbook-xml-to-rtf for more details).

7.7.1 SGML style sheets

CONFIRM THE BEHAVIOUR WRITTEN BELOW IS EXCACT

SGML style sheets are *written in DSSSL* (others?), which is a Lisp-dialect, and the usual Lisp-conventions apply. Several

several? *what can this do?*

The DocBook reference recommends using jade written by James Clark, to do DSSSL conversions. *long processing times later implemented in XT with threads to allow faster output*.

Show a sample stylesheet

7.7.2 XML style sheets

XSL - the XML style sheets - are written in XML too. The tags in the name space “xsl:” describe what is to be done with the source XML tree, in terms of tag-remapping (<para> to <p>), sorting of subtrees, and several programming constructions like “if”, “while” and “foreach”. New tags can be constructed, existing tags can be altered, and *yes?*

a reference listing here

This style sheet - lifted from the DocBook distribution - creates an outline of a DocBook document by looking for headline tags, and make a tree of them *Eh?* Line breaks have been introduced for readability.

```
<!--
| Identity Transform Stylesheet
| ~~~~~
| $Author: ravn $
| $Date: 2000/03/02 21:55:31 $
| $Source: /home/ravn/CVS/SPECIALE/SPECIALE/x-outline.xsl,v $
| $Revision: 1.1.1.1 $
+-->
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform"
                version='1.0'>

<xsl:output method="html"/>

<xsl:template match="/">
  <html>
    <body bgcolor="white">
      <xsl:apply-templates/>
```

```

    </body>
  </html>
</xsl:template>

<xsl:template
match="set|book|part|reference|article|sect1|sect2|sect3|sect4|sect5|section|
  simplesect|preface|chapter|appendix|example|figure|table|informaltable">
  <xsl:param name="treeicon">n</xsl:param>
  <xsl:variable name="level" select="count(ancestor-or-self::*)"/>
  <xsl:variable name="title" select="title"/>
  <xsl:variable name="this" select="local-name(.)"/>
  
  
  
  <xsl:variable name="color">
    <xsl:choose>
      <xsl:when test="$level = 1">navy</xsl:when>
      <xsl:when test="$level = 2">red</xsl:when>
      <xsl:when test="$level = 3">blue</xsl:when>
      <xsl:when test="$level = 4">black</xsl:when>
    </xsl:choose>
  </xsl:variable>
  <span style="color:{$color}"><xsl:value-of select="$title"/></span>
  <xsl:if test="local-name(.)='example' and programlisting/@role">
    <small style="font-family:courier"><xsl:value-of
      select="substring-after(programlisting/@role,'-')"/></small>
  </xsl:if>
  <br />
  <xsl:apply-templates/>
</xsl:template>

<xsl:template match="*|@*|comment()|processing-instruction()|text()"/>

</xsl:stylesheet>

```

Since the DocBook reference is very vague on this matter (the book went to press before the XSLT recommendation was finalized) the Internet has been a great help. Table 7.3 lists the tested XSLT processors which conforms to the W3C XSLT-19991102 recommendation.

XT	XSLT processor based on the XP parser. Both written in Java <i>by James Clark</i> .
Xalan	XSLT processor written in Java originally based on the LotusXSL processor by <i>alphaworks?</i> . Currently under heavy development by the Cocoon team (xml.apache.org).

Table 7.3: Leading XSLT processors implementing [XSLT-19991102](#) in February 2000

During testing it showed very quickly that the Java tools are quick enough to use for casual use (*do a few benchmarks*), even though they are not yet tuned for optimal performance. Additionally Even so, implementations in C should only be considered by designers if Java clearly cannot do

the job. The advantages of Java over C ([reference big list, probably at JavaSoft?](#)) in program development

- Very hard to do pointer access
- A large set of tested libraries

Xalan is at the time of this writing in a development state, and contains bugs which causes it to crash on some of my sample DocBook XML document. After that I turned to XT which is extremely stable and reliable, and that I have used for the rest of the project - my only complaint is that it stops after reporting the first error instead of parsing the whole document.

When the Cocoon project stabilizes they will have a high end XSLT-engine embedded as a servlet in Apache, which makes this a project to watch. The technology is at the time of this writing still immature, but very promising.

The platform independance of Java indirectly prompted the implementation of remote filters in Cactus. Very early in the evaluation Xalan threw the above mentioned exception when processing a medium sized DocBook document. In order to rule out errors in the underlying Java environment on Asserballe (see section A.3.1 on page 90 for technical details) an identical run was made with JDK 1.2 on Aalborg. The same error occurred, but in less than one second instead of ten. Experiments showed that Java programs consistently runs 15 times faster on Aalborg than on Asserballe.

[and so what?](#)

[incoming.xml](#)

```
<?xml version="1.0"?>
<?xml-stylesheet href="incoming.xml" type="text/xsl"?>

<?cocoon-process type="sql"?>
<?cocoon-process type="xslt"?>

<page>

  <connectiondefs>
    <connection name="foo_connection">
      <driver>org.gjt.mm.mysql.Driver</driver>
      <dburl>jdbc:mysql://localhost/Cactus</dburl>
      <username>XXXXXXXXXX</username>
      <password>XXXXXXXXXX</password>
    </connection>
  </connectiondefs>

  <query connection="foo_connection" tag-case="lower">
    select when, mime, user, how, filename, comment,length(item) as l from incoming
  </query>

Testing
</page>
```

incoming.xml

```
<?xml version="1.0" ?>
<xsl:stylesheet xmlns:xsl="http://www.w3.org/1999/XSL/Transform">
<?cocoon-process type="xslt"?>

<xsl:template match="page">
<html>
<head>
<title><xsl:value-of select="title"/></title>
</head>
<body bgcolor="#ffffff">

<xsl:apply-templates/>
</body>
</html>
</xsl:template>

<xsl:template match="ROWSET">
<table border="0" bgcolor="#F0F0F0">
<tr><th>when</th><th>mime</th><th>user</th><th>how</th><th>filename</th><th>comment</th><th></th>
<xsl:apply-templates/>
</table>
</xsl:template>

<xsl:template match="ROW">
<tr>
<td><xsl:value-of select="when"/></td>
<td><xsl:value-of select="mime"/></td>
<td><xsl:value-of select="user"/></td>
<td><xsl:value-of select="how"/></td>
<td><xsl:value-of select="filename"/></td>
<td><xsl:value-of select="comment"/></td>
<td><xsl:value-of select="1"/></td>
</tr>
</xsl:template>
</xsl:stylesheet>
```

7.8 Converting documents to XML

XML is currently having the same problems that HTML had in its initial years, namely lack of software support. The difference is that HTML is so lenient that it is easy to write HTML by hand in a text editor. Not so with XML.

The difference this time is that Microsoft has been an active participant in the *definition of XML*,

and has made [???] about Windows 2000 supporting XML directly *in what precisely?*

Even so there is still a severe lack of conversion tools for the many files out there in non-SGML based formats like Word, Excel, \LaTeX , \TeX , Lotus Notes, etc., which must be available before the documents can be converted.

Majix	Java based RTF to XML, which I have adapted to DocBook
pod2docbook	Converts Perl documentation to XML
wordview <i>others</i> <i>sp?</i>	Converts <i>Word 6-95-97</i> to HTML. I have adapted them to produce DocBook output <i>better do it too</i>
tex4h	package for \LaTeX which can <i>what was it now?</i>
<i>others</i>	??

Table 7.4: Conversion tools to XML [February 2000]

Due to the rather inflexible DTD-requirements in XML, it is most convenient that such a tool generates stand

Move filter descriptions to Cactus chapter discussing filters? ... wordperf 9 should be ablt ot do sgml

7.8.1 Majix - RTF/(DOC) to XML

Majix *URL, company description* is an RTF (for Word 97 documents) to XML converter. I have created a configuration and a small style-sheet which can render a RTF file to DocBook XML.

If run with the Microsoft Java Machine (“jview”) while logged into a Windows machine with access to the GUI, it can launch Word to convert a given DOC file to RTF and process it, with very reasonable results. I have experimented with making this accessible to Linux by installing Microsoft Services for Unix on a NT machine with Word installed, and using telnet to invoke Majix. Unfortunately this doesn’t work. Another approaches (*which will be tested if time allows*) includes having a pseudo-user watching a directory on a NT-server and launching Majix whenever DOC or RTF documents arrive. These can at best be described as kludges.

Even though it shows great promise, it has been a year without new releases (including the promised Pro version), and the source is not available.

Conclusion: The Majix software is currently best for personal end-user conversions.

7.8.2 pod2docbook - POD to XML

The POD format is created for writing documentation for the Perl programming language, and was created to be “... an idiot- proof common source for nroff, \TeX , and other markup languages, as used for online documentation” by Larry Wall (*the perlpod manual page*).

The original release produced SGML DocBook. I have submitted patches to the author which allows the user to choose between XML and SGML DocBook.

Conclusion: This do a good job for documents written in POD, and will probably with time take over as the default formatter for printed versions of Perl documentation.

7.8.3 tex4h - convert T_EX to XML

What daelen does it do

Conclusion:

7.8.4 wordview with friends

Parses binary stream and produces HTML, WML plus more.

7.9 Rendering HTML on the fly

look at docs.sun.com - what daelen do they do? Render SGML to HTML on the fly?

When the decision has been made to provide information in *ML on the web server, another decision must be made. Should the various renderings to PDF and HTML be created before they are requested, or when the user asks for a certain rendering?

Look for an article on static versus dynamic

As always it depends on the nature of the data.

Static renderings is probably suitable for documents which rarely change their rendering or take long time to render, and which just must be served as fast as possible. Often there is a desire for an inexpensive solution⁵, and in that case a standard Apache with a rich set of prerendered documents in a flat file system may very well be optimal. Personal experience has shown this to be a very robust solution which very rarely requires human intervention. esr Eric S. Raymond reports that a modern PC with

Dynamic rendering is suitable whenever the data changes very often, or there is a wish for the user to be able

The history and usage of SGML. Creation of XML. Describe document validation, and conversion (XSLT) to XML and other formats. Freedom from restrictions of HTML. Problems with tons and tons of DTD's. The need for well-documented standard, robust, supported DTD's (current ones: TEI, ebook [buh], DocBook). HTML conversion utilities can be tailored to generated DocBook XML for SSP (currently pod2docbook, Excel xls2xml). Found that SGML (jade) is powerful but too slow for on-the-fly stuff, as opposed to XML renderer. Several to choose from if they conform to the w3c standards (DOM and SAX). Who uses DocBook at the moment? Man pages in DocBook on Solaris (look on machine). IE50 cannot show "simple" DocBook XML yet but it is the goal of that project. the xsl/docbook/contrib/outline/outline.xsl is an excellent sample of a small, powerful style sheet.

⁵"Inexpensive" here is defined in terms of man hours needed to learn and maintain the web server, as well as the hardware needed

7.10 DocBook considerations

DocBook⁶ is a DTD designed for documenting software. Used by Sun, FreeBSD, O'Reilly and others.

The SGML is converted to HTML by applying a HTML style sheet for DocBook, and sending it through jade w

screen shot The generated navigation tags are a little troublesome in Lynx, but nice in Netscape. A “up”, “next” and “previous” are available combined with a “Top title”->“Chapter title” -> “Section title” navigation bar at the top.

The localization code is manually maintained.

It is possible to have a DocBook XML-file which can be processed both with XSL-stylesheets *and* DSSSL-stylesheets, by letting the !DOCTYPE header point to the DocBook XML DTD along with system-dependent path (which is required by XML). A sample header for DocBook is listed below:

```
<?xml version="1.0" encoding="ISO-8859-1" standalone="no"?>
<!DOCTYPE article
  PUBLIC "-//Norman Walsh//DTD Simplified DocBk XML V3.1.3.6//EN"
  "/home/ravn/sgml/docbk/db315/docbookx.dtd">
```

If the standalone="yes" field is set in the <?xml?>-tag instead, it means that the XML-tools will not validate the contents against the DTD **does this increase performance?**, but still apply the style sheets.

Format	Command
XHTML-1.0	
HTML-4.0	java com.jclark.xml.sax.Driver \$< ../docbook/html/docbook.xml \$@
XSL-FO	java com.jclark.xml.sax.Driver \$< ../docbook/fo/docbook.xml \$@
RTF	jade -t rtf -d \$(HOME)/sgml/dsssl/docbook/print/docbook.dsl \$(HOME)/jade/pubtext/xml.dcl \$<
MIF (Frame)	
TeX	
PDF	Two approaches possible
Others?	

Table 7.5: Output formats possible for a DocBook XML document

Generally XSL conversions (commands starting with java) are fast enough to be done “on-the-fly”, where DSSSL conversions (commands starting with jade) are **too slow for this (timing?)**.

7.11 Converting documents to other formats

7.11.1 SGML to XML

```
sx -biso-8859-1 -xlower input.sgml $>$ output.xml
```

The DTD for the document *must* be registered in the `catalog`. Even so a number of warnings and errors may be printed, since some SGML-constructions cannot be represented in XML, and some entities are unknown. This must be handled manually - SX is not a tool do this perfectly.

7.11.2 XML to RTF

Use jade with the appropriate stylesheet and `....`

```
jade -t rtf -d docbook.dsl xml.dcl inputfile.xml
```

where `docbook.dsl` is one of the stylesheets from the `DocBook Modular StyleSheets` (`dsssl/docbook/print/` and `xml.dcl` comes from the Jade distribution (`pubtext/xml.dcl`)

The resulting RTF-file is compatible with Word-97, Word-95 (claimed by author), StarOffice, `other texteditors`

`other backends`

SQL to XML :: Coocncococns sql proecessor

Chapter 8

Standard Query Language, Databases and Webservers

“You’ve got to be carefully taught”
South Pacific – Rogers and Hammerstein

8.1 SQL

SQL (the Standard Query Language) was developed in the early 1970’s by Don Chamberlin and Ray Boyce at IBM Research, to implement the mathematical notation in which E. F. Codd had defined relational databases [Codd(1970)]. Today SQL has grown into the *de-facto* standard for communicating with a relational database, which is being enforced with the modern ODBC and JDBC drivers which provides a “tunnel” between the client program and the database where the client can send a SQL-command and get the result back without caring about *how* the data are transmitted, much in the same way that the TCP/IP protocol allow one computer to send a data-stream to another computer without concering about how the data is sent.

Unfortunately this is as much abstraction as you will get from these “SQL drivers” – the ability to execute a SQL-command and retrieve the result. Any further abstraction from plain SQL is still left to the application programmer which means that even the state-of-the-art applications talk to the database in a command-line fashion. The SQL looks like this (lifted from the MySQL tutorial):

```
CREATE TABLE shop (  
  article INT(4) UNSIGNED ZEROFILL DEFAULT '0000' NOT NULL,  
  dealer CHAR(20) DEFAULT '' NOT NULL,  
  price DOUBLE(16,2) DEFAULT '0.00' NOT NULL,  
  PRIMARY KEY(article, dealer));  
  
INSERT INTO shop VALUES  
(1, 'A', 3.45), (1, 'B', 3.99), (2, 'A', 10.99), (3, 'B', 1.45), (3, 'C', 1.69),  
(3, 'D', 1.25), (4, 'D', 19.95);
```

```
SELECT * FROM shop;
```

article	dealer	price
0001	A	3.45
0001	B	3.99
0002	A	10.99
0003	B	1.45
0003	C	1.69
0003	D	1.25
0004	D	19.95

Figure 8.1: The output from the SQL-example

which then returns the table shown in figure 8.1:

The table shown here was originally in ASCII (manually put into a table for presentation purposes); inside an application program it is common for the application to work with a row at a time.

8.1.1 SQL - really a standard?

SQL has evolved over the years. The latest ANSI standard is SQL92, which is discussed in [Melton and Simon(1992), *as well as the standard bodies online*], and it provides several levels of compliancy a given SQL-database can adhere to. Even so, there are several facilities which are still not covered by the standard:

- **unique numbers** – having a row or function which generates a unique number every time it is called. This is called *sequence numbers* in Oracle, and `?` in MySQL.
- **date and time functions** – querying for the current date and handling date/time fields in the database is not standard. Additionally it is common for a vendor to have additional formats for timestamps which is incompatible with other vendors.
- `others?`

These vendor specific extensions from SQL92 are normally done in order to allow faster database queries, and are therefore interesting for the application programmer who wants the fastest possible implementation. Unfortunately this usually ties the application dependent to a specific database server, which rarely is in the interest of customers (who might have a database already they wish to use), meaning that the application programmer must be aware which facilities are acceptable to use. I have not found application frameworks which provide a generic SQL interface and parse it to the underlying database.

8.2 SQL – getting started

I was reintroduced to SQL rather abruptly in October 1999, where I needed one or more good introductory texts for ASP-programming with MS-Access as the database backend¹, and found almost none on the Internet. Martin Damsbo, who had previously written the PHP3+MySQL based WACO system [Damsbo(1999)], pointed me towards the tutorial in the MySQL reference manual², which was what he had used.

Even though several databases were available, I decided to delve into MySQL since it had excellent Perl support which I had chosen for developing the Cactus system (see 9.1.5 on page 80 for reasons), and was OpenSource software available on Linux.

The MySQL tutorial is reasonable for a person knowledgeable in SQL to get acquainted with the MySQL database, but it falls a bit short of being a generic introduction.

A *much* better choice for me was Phil Greenspun's SQL-introduction³ (see [Greenspun(1998)]) which was specifically written as introductory material for college students, based on his long-term real-life experience with the heavy duty web-service at ArsDigita⁴, which also serves photo.net⁵ where this SQL-introduction resides..

Greenspun starts with the basics and moves to advanced topics in a steady pace, and is an excellent read. There is a strong emphasis on Oracle 8i, but is easily adapted to other databases like MySQL. He has numerous examples, and several references and evaluations of other texts. Additionally the text is sprinkled with thumbnail photos, which links to a large version along with a link to another part of his site.

a swkveen shot?

This in addition with the comments from the readers to the text, makes this a very valuable resource for users. The SQL-introduction is available both online and as a book.

8.3 Some of the possible ways of accessing a database from a web-server

Armed with the necessary, solid SQL knowledge, it is reasonable to consider the current approaches of interfacing a Web server with a SQL-database. At least the following options exist:

ODBC – The Microsoft “Open Database Connectivity⁶” which provided for the first database access standard on Windows. Since superseded by OLE DB which does basically the same thing on Windows.

¹I learned quite a lot in those few days. Basically ASP is a very good idea, and Frontpage is extremely unsuited for working with it

²... the tutorial in the MySQL reference manual. ... – http://sunsite.auc.dk/mysql/Manual_chapter/manual_Tutorial.html#Tutorial

³... Phil Greenspun's SQL-introduction. ... – <http://photo.net/sql>

⁴... ArsDigita. ... – <http://www.arsdigita.com>

⁵... photo.net. ... – <http://photo.net>

⁶... Open Database Connectivity... – <http://www.microsoft.com/odbc>

JDBC – The “JDBC Data Access API”⁷ does the same thing for Java-programs as ODBC do for Windows-programs.

Perl DBI – The “Perl Database Interface” allows Perl programs to speak to Adabas, DBMaker, Fulcrum, Informix, Interbase, Oracle, Solid, Empress, Sybase, Illustra, Ovrmos, PostgreSQL, QBase.

Python DB-API – A rapidly growing set of database interface modules for Python⁹. Currently are Informix, Interbase, MySQL, Oracle, Sybase and the generic ODBC interface supported.

Vendor dependant – Each vendor have their own way of moving HTML response (and some XML) into the database itself, basically allowing the database server to respond to webrequests, but these are more often than not very slow.

8.4 What databases are available?

8.5 What database should we use?

Recommendations:

JDBC support is gaining widespread use in general, meaning that java programs can interface very easily.

If the implementation language and platform is fixed (like perl on solaris, or visual C on Windows) look at the driver support available and choose the DB from that. For development purposes the Java + MySQL comination has worked well for me.

⁷...JDBC Data Access API...–<http://www.javasoft.com/jdbc>

⁹...database interface modules for Python...–<http://www.python.org/topics/database/modules.html>

8.6 Can users transparently harness the power of XML for webpublishing?

Hmmm looks familiar. Integration wity DB instaead

Discuss how users can use their current tools to publish information to the web (Cactus principles: data capture and automatic conversion) by automatic conversion to XML and on-the-fly conversion on the way out of the webserver (while still having access to the original file provided by the user). This also allows people to view files made with software products not available to them.

8.7 Active and passive data acquisition for a database backed web-server

<p>[just thoughts] User can “push” data to the system, or the system can actively watch a resource (directory, webpage, usenet server, “channels”) and update whenever new data is available. www.mind-it.com. RDF-format from netscape to push headlines.</p>

8.8 Search engines and data mining

<p>htdig (process website) swish (process files for website) are ok. With database you can ask the database directly (currently in blobs) but with more analysis of the files you can do better. “Which author wrote this back in 1948”.</p>
--

Chapter 9

Sample implementation – Cactus

Dette afsnit er pt brutalt pillet ud af en Word fil jeg har skrevet tidligere, og lige justeret til. Indholdet haenger ikke sammen endnu. Bare lige for ikke at afskraekke :-)
--

squid: <http://www.squid-cache.org/>

9.1 abstract

The Cactus system fills a gap in the current integration of the web with normal office procedures, as it provides easy web publishing for occasional authors, by letting them submit documents in several ways with their usual software.

The system automatically produces other versions of the documents on demand of the users reading the documents, as well as provide the navigational framework, relieving the local webmaster from these tedious and error prone tasks.

9.1.1 System description

Cactus is a sample implementation of the following issues:

- **easy publishing** - users can publish via email/fax/print/watch usenet/www which is processed into the SQL database, and confirmed.
- **automatic conversion** - system will automatically convert a given document to what the user can see (or want). Browser sends a capability string with each request - use that to provide stuff directly, or generate a “this is available” summary.
- **automated navigational framework** - each document has an annotation which tells Cactus where to place it in the navigational hierarchy. The corresponding navigational pages are automatically generated and updated when new documents arrive. This also ensures system integrity without “broken links”.

Implementation languages - possibly Java or Perl. Perl chosen due to better library support (with source). Rex thingie in Java. How can MIME, TAR, etc be done in Perl (remember jublations!). Using Apache with Cocoon (perhaps) and MySQL. Graph algorithms in [Sedgewick(1990)].

Good summary from 19991114. Speed of PNG gzipping? PNG uncompressed initially and then later compressed by pngcrunch.

<http://photo.net/wtr/word.html> Utilities: mswordview, xls2xml

9.1.2 Background

Explain the history. Yggdrasil to address the need of an automatic webmaster → analysis (separate file written earlier). Cactus to address the difficulties of publishing information to Yggdrasil.

Explain Cactus as a successor to Yggdrasil.

9.1.3 Yggdrasil - a simple navigational framework

The Yggdrasil system was designed to provide a simple, consistent navigational framework around a dynamic set of web pages provided by users, as well as providing a “recently changed pages”-list plus an overview of pages written by each person. This was a very successful idea in the start, but gradually lost momentum due to these factors:

The development platform was changed from Unix to NT, which made it much more difficult for the users to access their web-directories, as these were no more a part of their home directory¹

The internal document format became Microsoft Word, which at that time could not be converted to HTML. Such documents were therefore unable to be published.

¹The Apache webserver uses the /public_html directory for the users personal web pages. The transition to NT meant that the users - in addition to their normal file manipulations - should telnet to a Unix server, and change the file attributes every time the file was updated.

The users could trigger an update by sending an empty email to the system, as well as rely on an automatic nightly update. The trigger mechanism was not brought along when the email system was converted to NT.

New users was not informed about the system.

9.1.4 Installation

Perlmodules. MySQL server. SAMBA.

[Rephrase following paragraph]

I have concluded the following goals for CACTUS, in order to avoid repeating history:

1. "Ease of publishing" is crucial.
2. Document preparations and transformations must be fully automatic.
3. The organisation using the system must be fully doing so.

"Ease of publishing" is crucial

CACTUS uses two approaches in order to make publishing as easy as possible for the users, namely

Documents are accepted in their native format, and in numerous ways.

Discussions with potential users showed [...] to be realistic ways in which CACTUS would be used:

As a document storage for various versions of a document, being developed and emailed back and forth between authors and peers. By allowing CACTUS to accept documents as email attachments, it would be very easy to enter each draft in CACTUS by including it on the list of authors or peers. In this way the archival of the document in CACTUS is completely transparent. As a fax machine. When replacing a fax machine with a computer, it is very easy to send a copy of the temporary image of the fax to CACTUS, before printing it. The fax system is then enhanced with the possibilities of the web, relieving the need for the physical copy. As a printer. Cactus provides a "printer", which makes a web-version out of any document the users can print. Users can then share final versions of documents without requiring the recipients to have the software in question. Either the Adobe Acrobat Reader can be used, or the primitive multiple image viewer in Cactus.

The full list of the publishing methods in Cactus is listed [....]

The users were primarily expecting to use these file formats:

- Word DOC and RTF files. These are the storage formats of the Microsoft word processor Word.
- HTML files. The common format for the web.
- LaTeX files. This typesetting program is very popular with mathematically oriented academics.

- GIF, JPEG, TIFF and PNG images. These and many more are in common use. Cactus use PNG internally [except possibly for JPEG]. The full list of supported formats is listed in the implementation section [see somewhere].
- PostScript and PDF files.
- Fax images.

Document preparation and conversion must be fully automatic

There was a strong consensus amongst the users, that it would be very nice not to have to convert the documents manually every time they were to publish a document. Therefore Cactus accepts several document formats as described above, and abandons the requirement that the users should convert their documents to HTML before publishing.

The system have some very different views on the documents depending on which representation the user needs - not all make sense for all documents:

1. The unaltered original. This file is always available, allowing users with the correct software to continue working with it.
2. A normalised version of a document. This could be a PNG version of a TIFF image or BMP image which can be viewed by all modern browsers, and an XML version of a document. If a suitable encoding can be found, even images and sound can be represented in XML so that this does not overlap the textual representation.
3. A visual representation of the original. This is the "look of the file", i.e. as it would look when printed to paper, and allows users without the corresponding software to view and print the contents.
4. A textual representation of the original. This is the "meaning of the file", which could say that the line "foo bar" is a level 2 heading, providing search capability.
5. An audio representation of a "document". This is e.g. an message left on an answering machine.

Conversions between all these document representation must be automated as much as possible. [write about the Cactus framework for providing existing and future filter types]. The normalised document is the only one created when a document enters Cactus - the rest are created on demand to avoid overfilling the underlying database, but cached for a reasonable time to improve performance.

[....]

The organisation using the system must be doing fully so.[rephrase]

In order for such a system to be successfully used within an organisation, it is vital that the organisation is using it whole-heartedly. [and more of the same].

Note on derivers:

A deriver is a pipeline which derives another version of an item. E.g. PNG to GIF,

There are the following tasks:

Derivers - derive another version of a given item. In theory, this is a reversible operation. These should generally be designed to be as fast as possible, since these will be called from CGI-programs with users waiting. Avoid compression (gzip should maximally be level 1).

Converters - create another item from one or more originals. This could be a multi-part MIME file which should be assembled into a fresh original. It could be a DVI file from several source files (tex files and images). These should generally be reasonably fast, since their speed specifies how fast a new item can be made available to the system.

Optimisers - A deriver is usually designed for being usable in a real-time interactive setting (also called being fast), which normally means that the result is sub-optimal. An optimiser complements this by doing a suitable optimisation step whenever the system is sufficiently idle. This could be running "pngcrunch" on PNG files (which yields 30% on Ghostscript output), or "tiffotiff" on TIFF files. The idea is that the file is still the same basic type, and can be used without further modification. It is an "in-place optimisation". Some file formats are using the gzip compression scheme internally. These would not benefit further from archiving.

Archivers - creates a long-term storage version of an item. This would typically be running "gzip -9" on the content, creating a new entry (with a new mimetype), and marking the item as detachable. The dearchiving process must be fast, since it might be needed by a deriver without notice.

Validators [spelling] - validates whether the content of the item is consistent with the MIME-type, and is it conforming to the standard. (Can a gzip-stream be decompressed? is a PostScript file parsing correctly? etc.). If a type is unknown, guess MIME-type from filename and/or contents, and validate it. If all fails, assign a type of application/octet-stream.

Examiner - looks into an item to look for references and encoded data. A reference may be an URL or an emailaddress in a signature. Encoded data could be an uuencoded image in a Usenet posting, or a pointer to a usenet article.

These are complemented by

Acquirers - gathers items from "outside" Cactus and enters it in the "incoming" table with an appropriate MIME-type. An acquirer could accept email, emulate a printer, retrieve Usenet articles, etc.

Janitors - cleans up whenever the system is otherwise idle, or when the cache is full. Derived items can be emptied of content, originals can be archived. Expired items can be purged completely from the database, along with all their derived items.

Presenters - extracts data, and present them. This could be a CGI-script presenting a given item as a http-stream. A PDF item could be presented as a window with two frames, the leftmost containing a thumbnail pr page, and the rightmost a high resolution version of a given page. The page should be cut in smaller pieces to allow easier processing by Unix Netscape.

In an ideal world we have infinite storage and infinite CPU-speed, meaning that everything would be done instantly when we need it. In order to decide the order in which to do tasks, a price system must be developed which would guarantee that the system processes every item, that the system is still responsive while converting,

What to do when the system runs so full that archiving cannot be done fast enough. Can data be moved to "outside storage?"

9.1.5 The implementation of Cactus.

Goals:

- Stable, well-known and remotely administrative platform- Linux/Solaris
- Platform independence - the resulting system must not be tied to Unix
- Modular - in order to minimise actual development, software libraries should be used as much as possible.
- Extensible - it should be easy for the system administrator to enhance functionality.

Choosing an implementation language:

Since platform independence was important for the system, it was quickly found that reasonable choices would include scripting languages plus Java. Scripting languages are fully interpreted allowing a script to run on numerous platforms without change. Java is the only compiled language with this property, due to the "Write once, Run everywhere" philosophy from Sun, plus the tight integration with Internet technologies in Java.

Initially I wanted to write the core of Cactus in Java, and spent a couple of weeks evaluating the language but found that

The mentality of the Java-community on the Internet, is very influenced by the shareware philosophy typical for the PC-user. Everything useful cost money, source code is not revealed, and the general tendency is for large, stand-alone applications.

Nobody had written a publicly available, stand-alone MIME-parser in Java. Several RFC's should be implemented and tested, in order to get email-parsing in Cactus.

The general level of abstraction is - in my opinion - too low in Java, while the "core language" is enormous.

[?]

Information extraction from items.

A given document contains one or more items, which again may contain further information. Cactus uses the MIME-type of a given item, to select the information scanning method with a fall back to the generic application/octet-stream examiner.

The application/octet-stream is parsed for:

The text is scanned for URL's, either with the Tom Christensen urlify or the program posted or commented by Abigail. These include ftp, http, gopher, mailto and news references. These are stored as external references.

The text/plain is parsed for

uuencoded data in a text stream. These start with ?begin ### name?, contain lines matching ####, and ends with ?end?. Such a file is then assigned a validated MIME-type based on the name of the file, and entered in the system as a derived item.

An text/html item is parsed for:

normal references in <a>-anchors. While doing this, the text to be rendered is extraced and parsed as a text/plain stream, in order to get the sequence right (may be changed). The <meta> tags are examined in order to extract keywords for the label.

9.1.6 Converters

Microsoft Word

The Microsoft Word DOC-format is widely used, but apparently so hard to use that even Microsoft cannot do it [reference to Word97 not being able to create Word95 files](#), and the reason why Cactus was started in the first place. I have spent a great deal of time looking for suitable software which could have been used with Cactus on the server side, and then testing it out.

Microsoft Word - The best solution would be able to actually run Word itself, but apparently this is integrated so well with Windows that it cannot run on other systems like WINE ([URL](#)). Since Microsoft previously have deliberately made the Windows 3.1 version of Internet Explorer 4 unable to run in the WinOS/2 subsystem for OS/2, I strongly suspect that this is also the case here.

Viewer for Microsoft Word - the Word Viewer is also available for 16 bit versions of Windows, which behave better in WINE. [A full Windows installation was not available to me when I tested this - t](#)
This can be used for producing the PostScript printouts.

Corel WordPerfect 8 for Linux - This is generally a very nice word processor, but the import filter for Word crashed WP when I tried it with a large document. Filters should be better in WordPerfect 9, which is due for Linux medio 2000.

StarOffice 5.1 - the German office suite for several platforms have excellent filters for importing Office files in general, but cannot be automated from the commandline. A StarBasic program must be written and invoked to do the job [ask on newsgroups.](#). This is currently untested, but probably the way to go. Since I did the testing, Sun have bought StarOffice - most likely since they need an network based office suite for Java, which StarOffice provide with thin clients and a Solaris based server - and promise to make the source generally available. This has not happened yet, but it is currently the most likely alternative to the Microsoft Office solution.

AbiWord - An Open Source word processor which shows great promise. It was not able to parse my test documents in Word.

mswordview [Others?](#) an Open Source Word to HTML converter which currently do text well, but have trouble with graphs which are converted to the WMF format.

Do-it-myself - If you ask Microsoft very nicely, you can get a copy of the Microsoft Office Binary File Format Specification, and write a personal parser of Word. It took Microsoft 6 months to

answer my email request, and then they just sent me the license twice with no specification. Second time they got it right. Unfortunately, the license was so restrictive that I decided not even to *look* at the specification.

Majix - This small RTF to XML converter can also parse DOC files, by invoking Word to save the DOC file as RTF. This shows great promise, except that this mechanism requires the GUI to be available for Word (it needs to pop up the window). The ingenious scenario of a telnet server on a NT-machine, which accepted remote logins from Cactus, started Majix which again started Word, did not work. Word hangs. Additionally the company has not released a new version in a year, and the license explicitly prohibits disassembly.

.... ..

My intermediate solution has been to adapt the mswordview to output DocBook XML instead of HTML, which provides a text-only display. Yes?

TeX and ~~TeX~~

9.1.7 MySQL

On asserballe

```
create table incoming (when datetime, mime varchar(80), user
    varchar(80), how varchar(80), filename varchar(255), comment varchar(80),
    item longblob);
```

Field	Type	Null	Key	Default	Extra
when	datetime	YES		NULL	
mime	varchar(80)	YES		NULL	
user	varchar(80)	YES		NULL	
how	varchar(80)	YES		NULL	
filename	varchar(255)	YES		NULL	
comment	varchar(80)	YES		NULL	
item	longblob	YES		NULL	

7 rows in set (0.01 sec)

- **wmf-anything:** libwmf (no fonts),
KVEC² (could not render the wmf files from wv),

²...KVEC... – <http://ourworld.compuserve.com/homepages/kkuhl/>

WMF docs³

9.2 Programs

9.2.1 Linux - an operating system

Linux was chosen as the development platform for these reasons:

- High degree of familiarity with the operating system
- Almost all Open Source software run “out-of-the-box” on Linux
- High performance Java Development Kit available from IBM
- Remotely maintainable via X/telnet/ssh.
- Freely available from the Internet.
- Root access possible without interfering with MIP security rules

These rules were also most likely fulfilled by any of FreeBSD/NetBSD/OpenBSD/Solaris x86, but I had already burned Redhat 6.1 on a CD for my portable, so there was no need to look any further. Linux has been a very satisfactory choice.

9.2.2 Apache - a high performance web server

There have probably been written more web servers than editors in the world of today. My requirements were simple:

- Run Perl programs efficiently
- Run Java servlets efficiently
- Have a large user base in order to ensure program availability

This basically left a single web-server, namely Apache⁴, which is capable of running Perl CGI scripts very efficiently with the `mod_perl` module⁵ (which uses a resident Perl interpreter combined with caching of the bytecode of previously encountered programs).

Other potential web servers? Roxen? Jetty?

9.2.3 Squid - a high performance web cache

The http-accelerator mode⁶

³...WMF docs... – <http://www.companionsoftware.com/PR/WMRC/WindowsMetafileFaq.html>

⁴...Apache... – <http://www.apache.org>

⁵...the `mod_perl` module... – <http://perl.apache.org/>

⁶...The http-accelerator mode... – <http://www.squid-cache.org/Doc/FAQ/FAQ-20.html#what-is-httpd-accelerator>

9.2.4 Cocoon

9.3 Filters

9.3.1 ps to pdf

Adobe Distiller – The reference PostScript to PDF converter. [myurl](#)

Ghostscript – This is a PostScript interpreter which has gained wide usage. [myurl](#)

rtf2latex <http://www.tex.ac.uk/tex-archive/support/rtf2latex2e/>

papirbjergene ImageTag, 3m xerox

Chapter 10

Conclusion

10.1 Conclusion

Databases are great, yea! Webservers are great, yea! Things must be powerful and easy for best results.

should this go away?

.1 \TeX and \LaTeX

which is well renowned for its excellence in general. For those who want a higher layer of abstraction from the raw, untamed \TeX the \LaTeX for

Appendix A

Report writing tools

When writing a large document like a thesis, you get to appreciate the power of the tools available to you in the Open Source community. This chapter outlines which tools I have used while writing and my experiences with them, with the thought that they might be useful to others later, and Open Source is all about not reinventing the wheel if it can be avoided.

All the software listed here, is available with a standard Linux Redhat 6.1 installation unless noted otherwise in the text.

A.1 \LaTeX

Guide is [Kopka and Daly(1999)].

A.1.1 Including screen dumps

Screen dumps have been captured with `xv`¹ (which is now 6 years old and still unbeat) and saved as `tiff` files. These were then converted to `png` and `eps` with the `NetPBM` package², which has worked flawlessly. It is important to remember to use `pnmtops -noturn rle` to get the correct orientation and a reasonable compression (25 files shrank from 52 Mb to 9 Mb total).

The `\includegraphic` (from the `graphicx` package) can read `png`-files from within `pdf \LaTeX` , and `eps`-files from within

The process of embedding Postscript is fully described in [Goosens et al.(1997)Goosens, Rahtz, and Mittelbach

I have had “`xdvi`” active while writing, and if the command “`latex mydocument && killall -USR1 xdvi`” is used, `xdvi` will automatically update after a successful compilation of the TeX-document. That is very nice when your editing have happened within a single physical page.

For presentation purposes, *Acroread with packages found at the DK-TUG website at aucdk*

¹...`xv`... – <http://www.freshmeat.net/appindex/1998/07/28/901622941.html>

²...`NetPBM` package... – <http://www.freshmeat.net/appindex/1999/06/24/930238715.html>

A.1.2 Finding back to the source document from xdvi

When reading the typeset document (either in “xdvi”, “acroread” or on paper) it is often a problem to find the appropriate tex-file. I used this trick which was deduced from the appropriate latex book here Daly? which puts the filename to the left in the running header on each page.

```
\usepackage{fancyhdr}           % provide headers where I can put the filenames.
% Trick from page 16 in fancyhdr documentation
\newcommand{\currentinputfile}{}
\newcommand{\myinclude}[1]{%
  \renewcommand{\currentinputfile}{File:\texttt{#1}}\include{#1}%
  \renewcommand{\currentinputfile}{}%

\lhead{\currentinputfile}
\pagestyle{fancy}

...
\begin{document}
\myinclude{s-front-matter}
\myinclude{s-introduction}
...
```

Modified xdvi!

A.2 Emacs

Emacs is a very well known editor in the Unix-world where its extensability allows it to get new functionality by installing additional modules. The standard distribution of Emacs provides:

Ispell – The M-Tab command (bound to TeX-complete-symbol) runs ispell-complete-word if the current word is not a TeX-command, which suggest possible words starting with the word under point. The M-\$ command (bound to ispell-word) is a great help when in doubt about a word. The flyspell-mode highlights words while you type, akin to the wavy underline used in MS-Word.

CVS – Emacs supports CVS³ directly. I had my working directory on my portable, and my CVS repository on the MIP system, in order for my CVS files to be automatically backed up.

Compile – The M-x compile works well with Jade, nsgml, Java and other tools where you may encounter errors in your source. Ctrl-x ‘ makes Emacs go to the line with the error.

Additionally I have installed:

AUC-TeX – provides a complete IDE for writing L^AT_EX documents, including keyboard shortcuts, a debugger, and intelligent commands depending on the state of the files. url? somewhere at sunsite. Highly recommended.

³ ...CVS...-???

psgmls – provides easy entering of tags, and parses DTD's to verify tags, and describe which tags are open at a given cursor location. Good for writing XML and SGML directly. The psgml homepage⁴ - *consider writing a set of AUC-TeX compatible keybindings for DocBook*

A.2.1 Making them all work together

Using these packages has basically been very easy, but occasionally they step on each others toes. That could normally be circumvented:

T_EX and RCS – The RCS header fields, contain dollar-signs which causes T_EX to enter/leave math-mode (like \$Id\$), causing it to be rendered in math italic *whichisnotsuitableforprose*. By typing it as \$ \$Id\$ \$ it is possible to avoid this. The dollar-space-dollar renders the space in math mode and switches back to normal text.

⁴... The psgml homepage... – http://www.lysator.liu.se/projects/about_psgml.html

A.3 Available hardware and software

This section lists the hardware I have used during development of Cactus.

A.3.1 Asserballe

Urls in long banes for this.

- Pentium-S 133 MHz PC with 128 Mb RAM and 10 Gb IDE disk
- Linux Redhat 6.1
- MySQL 2.3.?? - security updates
- Apache 1.?? with mod_perl, ApacheJServ, Cocoon.
- IBM JDK 1.1.8.
- CGI, DBI::??? perl modules for writing CGI scripts and accessing MySQL databases.

The choice of JDK was made at a time when top bug at Javasoft was “please provide a port of Java to Linux”, and where IBM beat Sun to it by porting their Java 1.1.6 implementation to Linux. The benchmarks showed it to be as fast as their Windows and OS/2 ports, which at that time was the fastest Just In Time compilers on the PC-platform. It performs very well on Linux, and the installed version is just a maintenance upgrade.

MySQL and the corresponding JDBC driver *must* be the latest stable versions due to security and bug fixes.

A.3.2 Aalborg

(4-CPU UltraSparc 440MHz with Sun JDK 1.2)

A.4 Bibliography with Web references

Full reference with a lot of annotated URLs.

Bibliography

- [Codd(1970)] E.F. Codd. A relational model of data for large shared data banks. <http://www.acm.org/classics/nov95/>, 1970. The paper that described relational databases mathematically. The history is explained in <http://www-4.ibm.com/software/data/pubs/papers/bontempo/>.
- [Damsbo(1999)] Martin Damsbo. Web access course organizer. <http://www.mip.sdu.dk/mez/waco>, 1999. MIP Intranet.
- [Dybvig(1996)] R. Kent Dybvig. *The Scheme programming language: ANSI Scheme*. Prentice-Hall PTR, Upper Saddle River, NJ 07458, USA, second edition, 1996.
- [Goosens et al.(1997)Goosens, Rahtz, and Mittelbach] Michel Goosens, Sebastian Rahtz, and Frank Mittelbach. *The L^AT_EX graphics companion - Illustrating Documents with T_EX and PostScript*. Addison-Wesley, 1997. Describes in great detail the intricacies of embedding graphics in LaTeX documents.
- [Greenspun(1998)] Phil Greenspun. Sql for web nerds. <http://photo.net/sql/>, 1998. Excellent introductory text to SQL.
- [Hofstadter(1979)] Douglas Hofstadter. *Gödel, Escher, Bach - an Eternal Golden Braid*, volume ISBN:?? ??, 1979. A very deep book which has influenced my thinking in many ways, both regarding computers, molecular biology and Bach, as well as the depth of self-reference. For fellow readers only: I am positive I have found the end of the book.
- [Knuth(1992)] Donald E. Knuth. *Literate Programming*, volume ISBN 0-937073-80-6. Stanford, California: Center for the Study of Language and Information, 1992. Knuths own page⁵.
- [Kopka and Daly(1999)] Helmut Kopka and Patrick W. Daly. *A guide to L^AT_EX*, volume ISBN: 0-201-39825-7. Addison-Wesley, 1999. An excellent book for the L^AT_EX author. I have used it a lot.
- [Melton and Simon(1997)] Jim Melton and Alan R. Simon. *Understanding the New SQL: A Complete Guide*. Number ISBN 1-55860-245-3. Morgan Kaufmann Publishers, San Francisco, California, 1997. A good introduction and reference to SQL-92. Defines the requirements for the various levels of SQL 92 compliancy .
- [Sedgewick(1990)] Robert Sedgewick. *Algorithms in C*. Addison-Wesley, 1990. ISBN 0-201-51425-7.

⁵...Knuths own page... – <http://www-cs-faculty.stanford.edu/knuth/lp.html>

- [Tanenbaum(1987)] Andrew S. Tanenbaum. *Operating systems - design and implementation*, volume ISBN: 0-13-637331-3. Prentice-Hall International, Inc., 1987. An excellent book in the inner workings of operating systems in general and Minix in particular. The complete source code is listed in the back of the book. *Minix was the direct inspiration of Linux.*

Contents

0.1	Abstract	2
1	Overview	3
1.1	What is a “database backed webserver”?	3
1.2	Why use a webserver in combination with a database?	3
1.3	The new role of the webserver	4
2	Terms and concepts	5
3	Dynamically generated documents	9
3.1	The way it started: Handmade web-pages in a filesystem on a web-server	9
3.2	When a site grows, it becomes hard to maintain	12
3.3	Navigational structure should not be maintained by the author	12
3.4	A good website needs meta-information about its documents	14
3.5	Avoid chaos by keeping information together	17
3.5.1	Javadoc - embedding web-information in programs	17
3.5.2	The Plain Old Documentation format for Perl	18
3.5.3	Literate Programming - Knuths approach to different views of the source . . .	19
3.5.4	Rational Rose - the other way around	20
3.5.5	Compiling source code into an executable	21
3.6	Requiring multiple views of a document	21
3.7	SGML/XML: Generic formats for storing data and meta-data	25
3.8	The user should use familiar tools to publish documents	26
4	The need for well-defined standards	27
4.1	Why is Open Source essential?	27

5	Sample websites	29
5.1	slashdot.org - high volume information site for nerds	29
5.2	Valueclick.com	30
5.3	Amazon - an Internet bookstore	31
5.4	Bang and Olufsen	33
5.5	Deja - where did Usenet go?	34
5.6	www.krak.dk	34
5.7	The Journey Planner - www.rejseplanen.dk	38
5.8	Freshmeat.net	40
5.9	Altavista - the search engine	41
5.10	The Collection of Computer Science Bibliographies	42
5.11	Jydske Bank	43
5.12	Cactus – document capture and conversion	43
6	Overview of technology available for Linux February 2000	45
6.1	Webservers	45
6.2	Database engines	45
6.2.1	Cloudscape - Informix	46
6.2.2	DB2 - IBM	46
6.2.3	Informix - Informix	46
6.2.4	Ingres - Ingres	46
6.2.5	mSQL - ?	46
6.2.6	MySQL - TCX	46
6.2.7	Oracle	47
6.2.8	SyBase - Sybase	47
6.2.9	Other approaches to database storage	47
6.3	Browsers	47
6.3.1	Netscape Communicator	47
6.3.2	Mozilla - Netscape 6.0	48
6.3.3	Opera	48
6.3.4	HotJava 3.0	48
6.3.5	Microsoft Internet Explorer 5.0	49

6.3.6	Amaya	50
6.3.7	StarOffice	50
6.3.8	W3 - browsing inside Emacs	50
6.3.9	Lynx - browsing in textmode	50
6.3.10	Perltools?	51
6.3.11	OTHER BROWSERS?	51
6.4	XML utilities	51
6.5	PDF utilities	51
7	SGML, XML and DTD's	53
7.1	"HTML is a SGML DTD" - the concepts	54
7.2	The author should provide content, not do layout	54
7.3	The creation of HTML	55
7.4	The creation of XML	56
7.5	XHTML - XML compliant HTML	56
7.6	Which dialect of SGML should be used?	57
7.7	How can a *ML document be viewed?	58
7.7.1	SGML style sheets	59
7.7.2	XML style sheets	59
7.8	Converting documents to XML	62
7.8.1	Majix - RTF/(DOC) to XML	63
7.8.2	pod2docbook - POD to XML	63
7.8.3	tex4h - convert T _E X to XML	64
7.8.4	wordview with friends	64
7.9	Rendering HTML on the fly	64
7.10	DocBook considerations	65
7.11	Converting documents to other formats	65
7.11.1	SGML to XML	65
7.11.2	XML to RTF	66

8	Standard Query Language, Databases and Webservers	67
8.1	SQL	67
8.1.1	SQL - really a standard?	68
8.2	SQL – getting started	69
8.3	Some of the possible ways of accessing a database from a webserver	69
8.4	What databases are available?	70
8.5	What database should we use?	70
8.6	Can users transparently harness the power of XML for webpublishing?	71
8.7	Active and passive data acquisition for a database backed webserver	72
8.8	Search engines and data mining	73
9	Sample implementation – Cactus	75
9.1	abstract	75
9.1.1	System description	76
9.1.2	Background	76
9.1.3	Yggdrasil - a simple navigational framework	76
9.1.4	Installation	77
9.1.5	The implementation of Cactus.	80
9.1.6	Converters	81
9.1.7	MySQL	82
9.2	Programs	83
9.2.1	Linux - an operating system	83
9.2.2	Apache - a high performance web server	83
9.2.3	Squid - a high performance web cache	83
9.2.4	Cocoon	84
9.3	Filters	84
9.3.1	ps to pdf	84
10	Conclusion	85
10.1	Conclusion	85
.1	\TeX and \LaTeX	86

A	Report writing tools	87
A.1	TEX	87
A.1.1	Including screen dumps	87
A.1.2	Finding back to the source document from xdvi	88
A.2	Emacs	88
A.2.1	Making them all work together	89
A.3	Available hardware and software	90
A.3.1	Asserballe	90
A.3.2	Aalborg	90
A.4	Bibliography with Web references	91
A.5	Cross references to this page have not been inserted at the correct page	99
A.6	The importance of a web cache	99

A.5 Cross references to this page have not been inserted at the correct page

A.6 The importance of a web cache

A database query is expensive, and it requires an expert to tune the database to run as fast as possible. It is not, however, always necessary to have the webserver do a database query to serve a page - often the generated page is valid for a short or long term period, and then it is relatively easy to cache the page for this period.

Here is one way to do it:

- Configure the script generating the page, to add an “Expires” header with a reasonable time of expiry
- Set up squid (9.2.3 on page 83) in http-accellerator mode, where it transparently adds cache facilities to a webserver, respecting the “Expires” header.

Mangler:

- myurl footnotes - urls should be broken if too long, numbers should perhaps not be reset at new chapters, slightly smaller bottom margin
- myimage maa gerne hyperlinke til tiff eller pngfil i hyperudgaven? ER det muligt?
- tal med sm om internetbibliotekarere
- afsnit om filosofi/osv med brugernes egne vaerktoejer til at lave sgml dokumenter.