

Findet Memo

Ein Archiv für Texte und andere Dokumente in Office-Formaten ist mit Oracle relativ einfach zu realisieren: Die Datenbank stellt dafür einen ganzen Werkzeugkasten mächtiger Funktionen bereit. Badran Farwati



Während Geschäftszahlen vielerorts seit Jahrzehnten mit Hilfe von Business-Software analysiert und verwaltet werden, blieben Texte lange außerhalb der Reichweite solcher Applikationen. Erst seit ein paar Jahren geraten zentrale, rechnergestützte Archive für Dokumente aller Art – Korrespondenz, Rechnungen, Berichte, E-Mails oder Notizen – stärker in den Blickpunkt. Schlagworte dafür sind etwa DMS (Document Management System) oder ECM (Enterprise Content Management). Derartige Anwendungen präsentieren ihre Dienste oft über Webportale und integrieren Workflow-Komponenten.

Wer auf solche Erweiterungen verzichten kann, außerdem bereits eine Oracle-Datenbank verwendet und ein wenig Kapazität für Skripting und SQL-Entwicklung frei hat, kommt aber schnell und kostengünstig auch mit einem Eigenbau zum Ziel. Denn Oracle bietet, was nicht jedem geläufig ist, ein Paket leistungsstarker Funktionen für die Suche nach

Texten in allen üblichen Office-Formaten – damit hat man das halbe Textarchiv schon in der Tasche.

Dieses Feature steht seit den frühen 8er Versionen auf der Liste, wurde allerdings von Version zu Version umbenannt: Erst hieß es Context (bis 8.1.5), dann Intermedia Text (bis 8.1.7), und schließlich Oracle Text for Oracle 8i, was sich inzwischen zur aktuellen Schreibweise Oracle Text verkürzt hat. Gemeint ist in jedem Fall ein Mechanismus, der das Indizieren von Texten und die Suche über diesen Index erlaubt. Die Metadaten der Texte – etwa Angaben zu Autor, Quelle, Titel, Entstehungszeit und so weiter – speichert die Datenbank in normalen Tabellenfeldern.

Quellenvielfalt

Je nach Größe finden die Texte selbst entweder in Textfeldern Platz oder werden in so genannten LOBs (Large Objects) abgelegt. Diese wiederum können

Bestandteil einer Tabelle sein (Datentyp BLOB, Binary Large Object; Maximalgröße 4 GByte) oder im Filesystem abgebildet werden. In diesem Fall enthält die Datenbank einen Verweis auf den Speicherort (Datentyp BFILE). Eine Datenquelle dieser Klasse nennt sich im Oracle-Text-Jargon Direct Datastore.

Alternativ können die Texte aber auch gewöhnliche Files (File Datastore) oder über das Inter- oder Intranet erreichbar sein (URL Datastore). Sie dürfen sogar auf Wunsch von einer Funktion geliefert werden, die der Anwender selber beisteuert (User Datastore).

Der Datenstrom, den Oracle aus einer dieser Quellen liest, wird dabei zunächst gefiltert. Passende Filter gibt es für mehr als 150 Datenformate, angefangen von Ascii-Text über PDF bis Word oder Excel. Zusätzliche Filter bieten mehrere Dritthersteller an. Den nächsten Arbeitsschritt übernimmt ein so genannter Sectioner, der in HTML- und XML-Dokumenten durch Tags strukturierte Abschnitte erkennt, nach denen zu suchen ebenfalls möglich ist.

Wörter, Themen, Klassen

Unstrukturierte Texte zerlegt der nachfolgende Lexer in kleine Einheiten (Tokens). Dabei berücksichtigt er die Eigenarten verschiedener Sprachen (bis hin zu Japanisch oder Koreanisch). Bestimmte Tokens lassen sich über konfigurierbare Stopp-Listen ausschließen. Am Ende stehen einzelne Wörter, aus denen ein inverser Index aufgebaut wird, wobei jedem Wort auch eine Liste seiner Fundstellen zugeordnet ist (siehe [Abbildung 1](#)).

Zusätzlich zur Volltextsuche bietet Oracle die Suche nach Themen an (per De-

fault allerdings nur in den Sprachen Englisch und Französisch). Dafür wird eine hierarchische Liste von Kategorien herangezogen, die es in Grenzen gestattet abzufragen, wovon der Text handelt. Ebenso möglich ist schließlich die direkte Klassifizierung von Texten. So könnten beispielsweise eingehende Nachrichten sortiert werden – eine Reihe spezieller Abfragen würde sie dann automatisch Ressorts wie Sport, Politik, Lokales oder Kultur zuordnen.

Volltextsuche

Das hier vorgestellte Beispiel behandelt nur einen Aspekt des komplexen Instrumentariums, die Volltextsuche. Da sich das Beispiel auf eine Prinziplösung mit exemplarischen Funktionen beschränkt, fehlt eine Reihe von Funktionen, die eine vollständige Applikation haben müsste. Beispielsweise ein Benutzer- und Berechtigungskonzept, eine Sessionverwaltung oder eine ausgefeilte grafische Oberfläche. Das alles und mehr kann aber entsprechend den eigenen Bedürfnisse ergänzt werden.

Das Beispiel stützt sich auf PHP für das Browser-GUI. Diese Skriptsprache verbinden die meisten sicherlich zuerst mit MySQL, sie bietet jedoch auch Schnittstellen zu einer ganzen Palette von Datenbanken. Der Oracle-Support wird beim Übersetzen von PHP mittels »./configure .. --with-oci8 = \$ORACLE_HOME« eingebunden [1].

Im vorgestellten Beispiel erfolgt die Ablage der Texte im Filesystem des Datenbank-Servers. Das Webinterface bietet zu diesem Zweck die Möglichkeit, Dateien dorthin zu importieren.

Tabellenbau

Zuerst müssen die zu archivierenden Dokumente auf den Server und in die Datenbank gelangen. Listing 1 stellt ein SQL-Skript vor, das zunächst die nötige Tabellenstruktur erzeugt. Die Tabelle »mein_doc« nimmt sowohl die Metadaten als auch die Verweise auf den Ablageort im Filesystem auf. Der wiederum ist in der Oracle-Prozedur »BFILENAME« als Parameter in Form eines Aliasnamens anzugeben, der sich mittels »CRE-

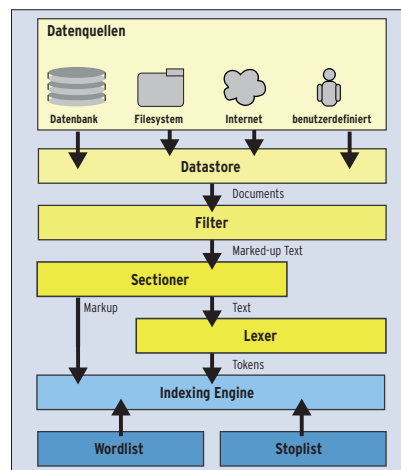


Abbildung 1: So entsteht ein Oracle-Index für Text. Die Daten wandern aus einer der möglichen Quellen durch eine konfigurierbare Verarbeitungskette und werden am Ende wortweise im Index erfasst.

ATE DIRECTORY« erzeugen lässt. Das Abspeichern übernimmt »put_file()«. Gelöschte Datensätze werden vorsichtshalber in die Tabelle »mein_doc_h« kopiert, wo sie nach Ablauf einer Sicherheitsfrist entsorgt werden können. Unter [2] finden sich auch alternative Skripte, die demonstrieren, wie Tabellen mit

Listing 1: »bfile_loesung.sql«

```

01 DROP TABLE mein_doc;
02 DROP TABLE mein_doc_h;
03 DROP SEQUENCE mein_doc_seq;
04 DROP INDEX mein_doc_idx;
05 CREATE TABLE mein_doc (
06   id NUMBER,
07   file_name VARCHAR2(21),
08   user_file_name VARCHAR2(255),
09   upload_datum VARCHAR2(10),
10   titel VARCHAR2(255),
11   autor VARCHAR2(255),
12   beschreibung VARCHAR2(255),
13   groesse VARCHAR2(20),
14   mime VARCHAR2(50),
15   inhalt BFILE,
16   CONSTRAINT doc_pk PRIMARY KEY (id)
17 );
18
19 create table mein_doc_h (
20   id NUMBER,
21   file_name VARCHAR2(21),
22   user_file_name VARCHAR2(255),
23   upload_datum VARCHAR2(10),
24   titel VARCHAR2(255),
25   autor VARCHAR2(255),
26   beschreibung VARCHAR2(255),
27   groesse VARCHAR2(20),
28   mime VARCHAR2(50)
29 );
30
31 CREATE SEQUENCE mein_doc_seq;
32
33 CREATE OR REPLACE DIRECTORY dokumente AS
34   '&1';
35 GRANT READ ON DIRECTORY dokumente TO ctxsys;
36
37 CREATE INDEX mein_doc_idx ON mein_doc(inhalt)
38   INDEXTYPE IS CTXSYS.CONTEXT
39   PARAMETERS ('SYNC ( ON COMMIT)');
40 CREATE OR REPLACE PROCEDURE put_file
41 (
42   p_file_name IN mein_doc.file_name%TYPE,
43   p_user_file_name IN
44     mein_doc.user_file_name%TYPE,
45   p_upload_datum IN
46     mein_doc.upload_datum%TYPE,
47   p_titel IN mein_doc.titel%TYPE,
48   p_autor IN mein_doc.autor%TYPE,
49   p_beschreibung IN
50     mein_doc.beschreibung%TYPE,
51   p_groesse IN mein_doc.groesse%TYPE,
52   p_mime IN mein_doc.mime%TYPE
53 ) AS
54 BEGIN
55   INSERT INTO mein_doc (id, file_name,
56     user_file_name, upload_datum,titel,
57     autor, beschreibung,groesse, mime, inhalt)
58   VALUES (mein_doc_seq.NEXTVAL, p_file_name,
59     p_user_file_name, p_upload_datum,
60     p_titel , p_autor, p_beschreibung,
61     p_groesse, p_mime,
62     BFILENAME('DOKUMENTE',p_file_name));
63 COMMIT;
64
65 -- Für Oracle < 10 aktivieren:
66 --
67   CTX_DDL.SYNC_INDEX(idx_name=>index_name);
68 END;
69 /
70
71 CREATE OR REPLACE PROCEDURE history
72 (
73   p_id IN mein_doc.id%TYPE
74 ) AS
75 BEGIN
76   insert into mein_doc_h select id, file_name,
77     user_file_name, upload_datum,
78     titel, autor, beschreibung, groesse, mime
79   from mein_doc where id = p_id;
80   delete from mein_doc where id = p_id;
81   commit;
82 END;
83 /

```


und Präpositionen, die allein keinen inhaltlichen Sinn haben – von der Indizierung ausschließen.

Wird der Index häufig synchronisiert, kommt es zu einer Fragmentierung. Deren Grad kann man mit Hilfe der folgenden SQL-Funktion ermitteln (Listing 2). Als Gegenmaßnahme bei zu starker Fragmentierung liefert Oracle die Prozedur »CTX_DDL.OPTIMIZE_INDEX« mit. Mögliche Parameter sind »FAST« oder »FULL«, zum Beispiel »EXEC CTX_DDL.OPTIMIZE_INDEX(idx_name => 'anno_doc_idx', optlevel => 'FULL');«.

Schließlich ist es empfehlenswert, von Zeit zu Zeit zu kontrollieren, ob es beim Indizieren zu Fehlern gekommen ist. Dies verrät ein Blick in die Tabelle »CTX_USER_INDEX_ERRORS«:

```
SELECT err_index_name, err_timestamp,
err_text
FROM ctx_user_index_errors
ORDER BY err_timestamp;
```

Ein Index für jeden Zweck

Der vorgestellte Index vom Typ Context beherrscht außerdem die Suche nach Wörtern des gleichen Wortstamms oder nach ähnlich klingenden Wörtern (so genannte Fuzzy-Suche). Er kann auch größere Datentypen bis zu 4 GByte indizieren und eignet sich daher für

»CLOB«, »BLOB«, »BFILE«, »VARCHAR2«, »CHAR« und für »XMLTYPE«.

Daneben kennt Oracle eine Reihe weiterer, spezialisierter Indexarten. Ctxcat eignet sich für die Indizierung von kürzeren Texten bei Abfragen, in denen Texte und Metadaten gemischt werden. In diesen Fällen verspricht er eine bessere Performance als ein Context-Index. Dafür verbraucht er aber 10- bis 15-mal mehr Platz als die Tabelle selbst und kennt weniger Suchoptionen als Context. Context benötigt dagegen maximal so viel Platz wie die Tabelle.

Fragewörter

Der dritte Indextyp, Ctxrule, kann nur auf »CLOB«- und »VARCHAR2«-Datentypen aufgebaut werden. Dieser Indextyp ist für die Klassifikation von Dokumenten bestimmt. Außer der Tabelle mit den Dokumenten legt Oracle eine zweite Tabelle an, in der die verschiedenen Dokumentklassen und vordefinierten Queries gespeichert sind, deren Ergebnisse der Zuordnung eines Dokuments zu einer der Klassen dienen.

Für die Volltextsuche kommen die üblichen Operatoren wie »LIKE« oder »=« nicht in Frage. Stattdessen verlangt jeder Indextyp einen eigenen Suchoperator: Ctxcat erwartet in der »WHERE«-Klausel den Operator »CATSEARCH«, Context

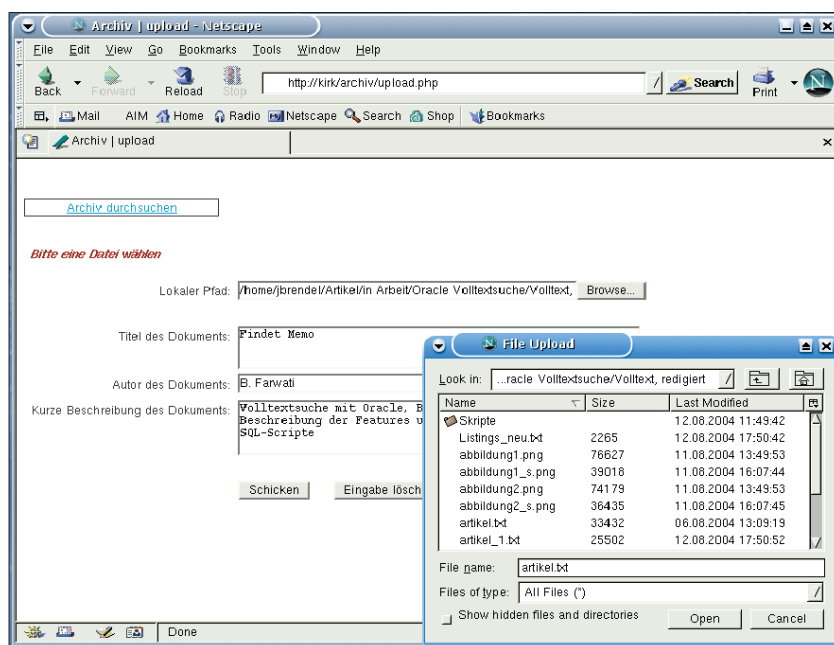


Abbildung 2: Über dieses Upload-Interface kann der Benutzer eigene Dokumente in das Archiv aufnehmen, in dem er sie später mit der Oracle-Volltextsuche durchforsten kann.

wird mit Hilfe von »CONTAINS« befragt und Ctxrule braucht »MATCHES«. Daneben gibt es »ABOUT« für die Themensuche und »WITHIN« für die Suche nach Abschnitten von HTML- oder XML-Dokumenten.

Datentypen wie »CLOB« oder »BLOB« lassen sich nicht in einer Ergebnistabelle darstellen, sie dürfen daher auch nicht als gewünschte Ergebnis-Spalten hinter dem »SELECT« vorkommen. Oracle wertet die Abfrage zuerst nach der Rangfolge der Operatoren aus. Gleichwertige Operatoren werden von links nach rechts gelesen.

Die wichtigsten Operatoren sind das logische »OR« (|), »AND« (&), »ACCUMulate« (,), »\$« als Symbol für die Stammsuchung, »NOT« mit dem Symbol »~« und »!« für die Suche nach ähnlichen Wörtern (Soundex).

Rangordnung

Oracle unterteilt die Operatoren in zwei Gruppen: Jene, die nur einen Operanden verlangen, und jene, die mindestens zwei brauchen. Innerhalb der Operatorgruppen gibt es jeweils eindeutige

Prioritätsregeln: Für unäre Operatoren gilt Stem vor Soundex; für die andere Gruppe gilt NOT vor AND vor OR vor ACCU [4].

Einige Beispiele demonstrieren die Verwendung dieser Operatoren. Eine bestimmte Zeichenkette findet man so:

```
SELECT id, file_name
FROM tabellen_name
WHERE CONTAINS(spalten_name,
'Universität Wien' ) > 0
ORDER BY id;
```

Wer nach mehreren Wörtern zugleich fahnden will, verknüpft sie mit dem Operator »&«:

```
SELECT id, file_name
FROM tabellen_name
WHERE CONTAINS(spalten_name,
'Universität & Wien' ) > 0
ORDER BY id;
```

Beim Testen dieser Abfrage mit SQLPlus ist eine wichtige Kleinigkeit zu beachten: Für SQLPlus ist »&« der Anfang eines Variablenamens, weshalb hier der vorangestellte Befehl »scan off« notwendig ist. Reservierte Wörter gehören in geschwungene Klammern.

Oracle kann die Fundstellen nach ihrer Relevanz bewerten, indem es einen

Punktwert (Score) pro Suchresultat berechnet. Der Programmierer bindet die Trefferbewertung an ein so genanntes Label – im folgenden Beispiel »1«. Eine entsprechende Abfrage sieht so aus:

```
SELECT SCORE(1), title
FROM tabellen_name
WHERE CONTAINS(spalten_name,
'oracle', 1) > 0
ORDER BY SCORE(1) DESC;
```

In Listing »get.php« (zu finden unter [2]) ist ein bescheidenes, allgemeines Webfrontend skizziert, das die Funktionen Stem, Soundex und AND verwendet. Diese Suchmöglichkeiten kann jeder nach seinen Erfordernissen weiter ausbauen. Ein Beispiel dafür wäre:

```
SELECT id, file_name
FROM mein_doc
WHERE CONTAINS(inhalt, '(urlaub & reise)
| ( hotel & appartement ) ~schnee') > 0 ;
```

Select findet Dokumente, die entweder die zwei Wörter Urlaub und Reisen beziehungsweise Hotel und Appartement (oder alle vier) enthalten, nicht aber das Wort Schnee.

Wer das hier beschriebene Beispiel auf seinem Rechner installieren möchte, ko-

Listing 4: »class_db.inc«

```
01 <?
02 include_once('./config.inc');
03 // Database Class for oracle
04 class db {
05     var $db;
06     var $user;
07     var $pass;
08     var $error_msg;
09
10     function db(){
11         global $DB, $USER, $PASS;
12         $this->db = $DB;
13         $this->user = $USER;
14         $this->pass = $PASS;
15         $this->open();
16     }
17     function open(){
18         $conn = OCILogon($this->user,$this->pass,$this->db);
19         $this->conn = $conn;
20         if ($conn == false) {
21             echo OCIErr($conn);
22             return false;
23         }else return true;
24     }
25     function query($query) {
26         $this->query = $query;
27         $handler = OCIParse($this->conn,$query);
28         OCIExecute($handler,OCI_DEFAULT);
29         $i =0;
30         while(OCIFetchInto($handler, $result, OCI_RETURN_NULLS +
31 OCI_ASSOC)) {
32             foreach($result as $key=>$value){
33                 $k = strtolower($key);
34                 $ergebnis[$i][$k] = $value;
35             }
36             $i++;
37         }
38         return ($ergebnis);
39         OCIFreeStatement($handler);
40         OCILogoff($this->conn);
41     }
42     function insert($insert_stmt){
43         $handler = OCIParse($this->conn,$insert_stmt);
44         //OCIExecute($handler,OCI_DEFAULT);
45         if(!is_resource($handler) || !OCIExecute($handler,
OCI_DEFAULT))
46 {
47             $this->error_msg = "SQL-Statement-Handler ist nicht
48 korrekt";
49             echo $this->error_msg;
50             return false;
51         }else return true;
52         OCIFreeStatement($handler);
53         OCILogoff($this->conn);
54     }
55 }
56 ?>
```

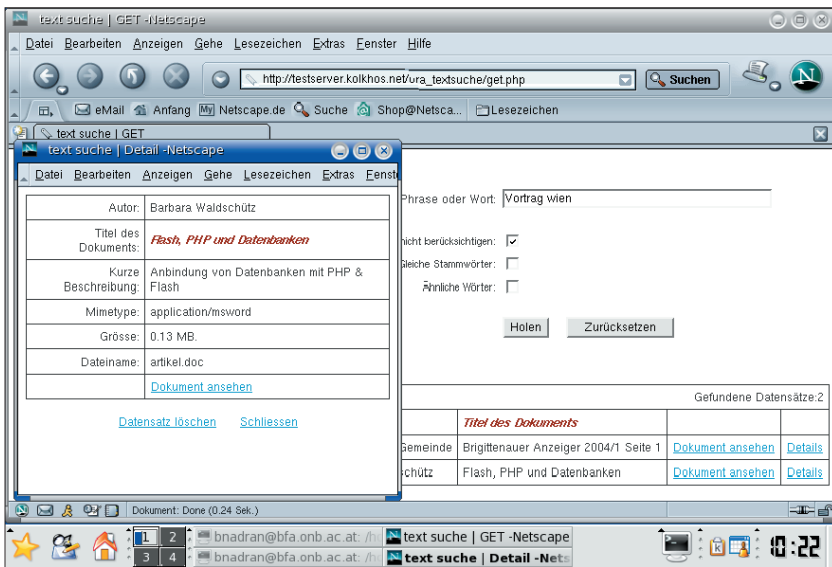



Abbildung 3: Die Suchmaske der Beispiel-Applikation: Rechts die Eingabefelder für die Suchbegriffe, links eine Trefferanzeige, die auch Metadaten des gefundenen Dokuments zeigt.

piert sich zuerst alle ».php«- und ».inc«-Dateien aus [2] in ein Unterverzeichnis des »htdocs«-Directory des Webservers und loggt sich dann als Oracle-User »system« in die Datenbank ein.

Anpassen, loslegen

Mit Hilfe des SQL-Skripts »create_user.sql« [2] legt man einen neuen User für die Applikation an. Dieses Skript erfordert vier Parameter: »@pfad/create_user.sql Username Passwort default_tablespace temp_tablespace«. Anschließend führt der neu angelegte User das Skript »bfile_loesung.sql« aus, das als Parameter den Pfad zu jenem Verzeichnis verlangt, in dem die Dokumente abgelegt werden sollen: »@bfile_loesung.sql Pfad zum Upload-Verzeichnis«.

Nun sind noch das Upload-Verzeichnis selbst und darunter ein Unterverzeichnis zu erzeugen, in dem später die aus der Datenbank gelöschten Dokumente landen sollen (im Beispiel das Verzeichnis »deleted«). Das Upload-Verzeichnis muss für den User beschreibbar sein, den den Webserver gestartet hat:

```
mkdir /upload
mkdir /upload/deleted
chown -R nobody:nobody
```

Darauf folgt die Anpassung der »conf.inc«-Datei. Hier gehören vor allem die jeweilige SID der Datenbank sowie Username und Passwort des Beispielbenut-

zers hinein. Zum Schluss legt man in der Datei »httpd.conf« des Apache-Servers noch einen Alias für das Upload-Verzeichnis an:

```
Alias /upload/ "pfad/upload/verzeichnis"
<Directory "pfad/upload/verzeichnis">
  Options Indexes MultiViews
  AllowOverride None
  Order allow,deny
  Allow from all
</Directory>
```

Die Änderungen werden nach einem Neustart des Webservers wirksam.

Buchstabensalat

Bei der Oracle-Installation – sie ist hier nicht näher beschrieben – sollte der Anwender darauf achten, dass er die Datenbank am besten mit dem Character-Set UTF-8 (Option »AL32UTF8«) anlegt. Dann kann die Applikationen Dokumente in unterschiedlichen Sprachen archivieren und durchsuchen. Nachträglich lässt sich das Character-Set der Datenbank nur dann ändern, wenn das alte Character-Set eine Untermenge des neuen ist, wobei sich alle gemeinsamen Zeichen an identischen Positionen befinden müssen.

Ein Wechsel etwa von ISO88591 zu UTF-8 ist deshalb nach der Installation der Datenbank nicht mehr möglich. Doch nicht nur bei der Datenbank ist auf den Zeichensatz zu achten: Auch Webserver und PHP müssen sowohl beim Doku-

menten-Upload als auch während der Suche UTF-8 verwenden.

Zwar werden im Beispiel alle Dokumente beim Upload mit »charset = utf-8« im HTML-Metatag gekennzeichnet, doch würde das überhaupt nichts bewirken, wenn in der »php.ini« gleichzeitig »default_charset = "iso-8859-1"« vorgegeben wäre. Denn seit der Version 4.0b4 ersetzt PHP automatisch die Zeichensatz-Einstellung im HTML-Header durch seine eigene Vorgabe.

Ein Ausweg besteht darin, in der PHP-Konfigurationsdatei »php.ini« die Definition des Standard-Zeichensatz frei zu lassen (genau ein Leer- zwischen zwei Anführungszeichen setzen). Das bewirkt, dass PHP auf eigene Eintragungen verzichtet und so die in den HTML-Seiten bereits vorhandenen Kennzeichnungen erhalten bleiben.

Im Unterschied zur Beispielapplikation wären eigentlich Datenbankverbindungen wünschenswert, die nicht für jede Abfrage neu geöffnet werden müssen,

sondern so lange offen bleiben, bis der Benutzer sich ausloggt oder bei Inaktivität ein sinnvoller Timeout-Wert überschritten wird.

Sitzungsstress

Diese Technik ist allerdings nicht besonders einfach zu realisieren. Zum Beispiel müsste die maximale Anzahl der Datenbank-Verbindungen abhängig vom verfügbaren Hauptspeicher des Servers berechnet und eine Reihe von Oracle-Parametern daran angepasst werden. Der Webserver wäre ebenfalls gezwungen, sich an dieses vorgegebene Limit zu halten. Möglich ist auch ein Connection-Pool, aus dem sich die Client-Prozesse bedienen können.

Oracle empfiehlt, den Datenbankserver in einem derartigen Szenario generell als so genannten Multi-Threaded-Server (MTS) zu konfigurieren. Weitergehende Hinweise dazu finden sich beispielsweise auf [6]. (jcb) ■

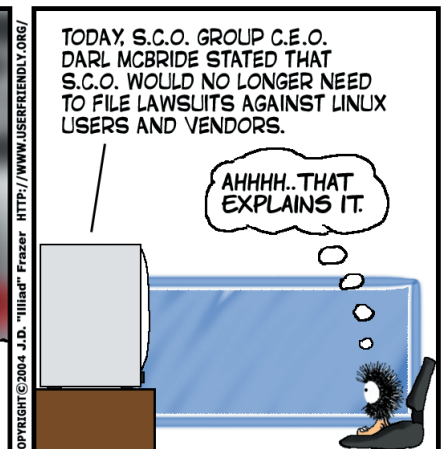
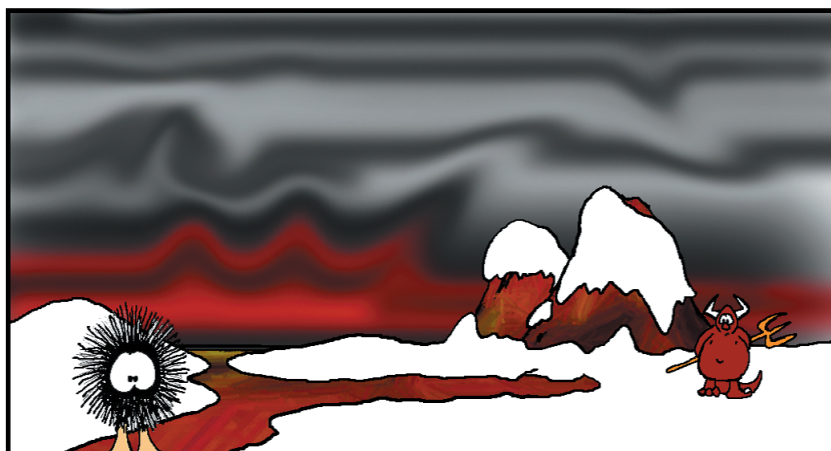
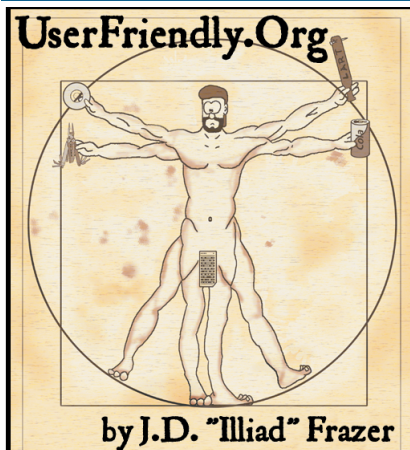
Infos

- [1] PHP-Manual: [<http://www.php.net/manual/de/ref.oracle.php>]
- [2] Listings zum Artikel: [<ftp://ftp.linux-magazin.de/pub/listings/magazin/2004/10/Volltextsuche>]
- [3] Oracle Text Application Developer's Guide: [http://download-west.oracle.com/docs/cd/B13789_01/text.101/b10729/toc.htm]
- [4] Funktionen des CTX_DDL-Package: [<http://www.cise.ufl.edu/help/database/oracle-docs/text.920/a96518/cddl/pkg.htm#11769>]
- [5] Oracle Text Reference 10g: [<https://cwisdb.cc.kuleuven.ac.be/ora10doc/text.101/b10730/title.htm>]
- [6] Oracle Multi-Threaded-Server: [<http://www.oracleadvice.com/Tips/MTS.htm>]

Der Autor

Mag. Badran Farwati arbeitet als Datenbankadministrator und -Programmierer in der Österreichischen Nationalbibliothek und beteiligt sich in seiner Freizeit an mehreren Open-Source-Projekten.

User Friendly - der monatliche Comic im Linux-Magazin



COPYRIGHT © 2004 J.D. "Illiad" Frazer. HTTP://WWW.USERFRIENDLY.ORG/