

Database Concepts for a new generation of CMS

Drupal - eZPublish - ...

Thomas Koch

Young Media Concepts GmbH

July 31, 2008

IT mögen.



- 1 Intro
- 2 Static database schemes
- 3 On Flexibility
- 4 Linking objects
- 5 Keys

Outline

- 1 **Intro**
- 2 Static database schemes
- 3 On Flexibility
- 4 Linking objects
- 5 Keys

What this talk is about

How to use a relational database effectively for a content management system.

What this talk is not about

- Alternative storage systems (CouchDB, XML)
- Other components of a CMS then storage

What does a CMS?

- Store content (Content = highly structured, mostly textual informations)
- Organize content
- Organize content relations
- Present content in different ways
- Implement Workflows, Permissions
- Search content

Examples of CMS

- CMS: Drupal, eZPublish, Plone
- Groupware: eGroupWare, Horde, Exchange
- Issue Tracker: Flyspray, OTRS, Bugzilla
- ERP: Organize documents: and workflows

What are not CMS's?

- OpenOffice, GIMP: Editors
- Gallery2: Organizes Images, not textual content
- Mailreaders, RSS reader, Browser: retrieve and present
- Not many Webapplications that are not CMS!
- more ... ?

Outline

- 1 Intro
- 2 Static database schemes**
- 3 On Flexibility
- 4 Linking objects
- 5 Keys

Tables are static

- Drupal: node, node_revisions, comments, poll, book, ...
- eGroupWare: egw_addressbook, egw_addressbook_extra, egw_cal, egw_cal_extra, ...
- OTRS: ticket, service, faq_item, ...

New modules come with new table definitions.

Schema for a Drupal module

```
function subscriptions_schema() {  
  $schema['subscriptions'] = array(  
    'fields'      => array(  
      'sid'       => array('type' => 'serial', 'unsigned' => TRUE,  
                          'not null' => TRUE, 'disp-width' => '10'),  
      'module'   => array('type' => 'varchar', 'length' => '64',  
                          'not null' => FALSE),  
      'field'    => array('type' => 'varchar', 'length' => '32',  
                          'not null' => FALSE),  
      'value'    => array('type' => 'varchar', 'length' => '237',  
                          'not null' => FALSE),  
      ...more fields  
    ),  
    'primary key' => array('sid'),  
    'indexes'     => array(  
      'module' => array('module', 'field', 'value'),  
      'recipient_uid' => array('recipient_uid')),  
  );  
}
```

Schema for an eGroupWare module

```
$phpgw_baseline = array(  
    'egw_bookmarks' => array(  
        'fd' => array(  
            'bm_id' => array('type' => 'auto', 'nullable' => False),  
            'bm_owner' => array('type' => 'int', 'precision' => 4, 'nullable' => True),  
            'bm_access' => array('type' => 'varchar', 'precision' => 255, 'nullable' => True),  
            'bm_url' => array('type' => 'varchar', 'precision' => 255, 'nullable' => True),  
            'bm_name' => array('type' => 'varchar', 'precision' => 255, 'nullable' => True),  
            'bm_desc' => array('type' => 'text', 'nullable' => True),  
            'bm_keywords' => array('type' => 'varchar', 'precision' => 255, 'nullable' => True),  
            'bm_category' => array('type' => 'int', 'precision' => 4, 'nullable' => True),  
            'bm_rating' => array('type' => 'int', 'precision' => 4, 'nullable' => True),  
            'bm_info' => array('type' => 'varchar', 'precision' => 255, 'nullable' => True),  
            'bm_visits' => array('type' => 'int', 'precision' => 4, 'nullable' => True)  
        ),  
        'pk' => array('bm_id'),  
        'fk' => array(),  
        'ix' => array(),  
        'uc' => array()  
    )  
);
```

conclusion

Modifying the database scheme in such systems requires programming and is a major surgery.

Outline

- 1 Intro
- 2 Static database schemes
- 3 On Flexibility**
- 4 Linking objects
- 5 Keys

columns in egw_addressbook

• contact_id	• adr_one_locality	• tel_prefer
• contact_tid	• adr_one_region	• contact_email
• contact_owner	• adr_one_postalcode	• contact_email_home
• contact_private	• adr_one_countryname	• contact_url
• cat_id	• contact_label	• contact_url_home
• n_family	• adr_two_street	• contact_freebusy_uri
• n_given	• adr_two_street2	• contact_calendar_uri
• n_middle	• adr_two_locality	• contact_note
• n_prefix	• adr_two_region	• contact_tz
• n_suffix	• adr_two_postalcode	• contact_geo
• n_fn	• adr_two_countryname	• contact_pubkey
• n_fileas	• tel_work	• contact_created
• contact_bday	• tel_cell	• contact_creator
• org_name	• tel_fax	• contact_modified
• org_unit	• tel_assistent	• contact_modifier
• contact_title	• tel_car	• contact_jpegphoto
• contact_role	• tel_pager	• account_id
• contact_assistent	• tel_home	• contact_etag
• contact_room	• tel_fax_home	
• adr_one_street	• tel_cell_private	

problems with static tables

- only two web addresses possible
- who needs tel_car and tel_pager anymore?
- No field for a client id

custom fields in eGroupWare

Table egw_addressbook_extra

- contact_id
- contact_owner
- contact_name
- contact_value

This is a first step in the direction of the *Entity-Attribute-Value* Pattern in database design.^a

^ahttp://en.wikipedia.org/wiki/Entity-Attribute-Value_model

Entity-Attribute-Value Pattern

Entity	Attribute	Value
1	n_family	Koch
1	n_given	Thomas
2	n_family	Musterman
2	n_given	Max
2	n_prefix	Dr.

eZ Publish uses this pattern in a more advanced version.

EAV Pattern - Advantages

- Flexibility

EAV Pattern - Disadvantages

- Many slow JOINS necessary
- Everything is a varchar
- Specialised tools necessary to browse the DB
- Hand coding SQL for this pattern is a pain
- Many build-in database features aren't available
 - referential integrity checks
 - domain integrity checks
 - indexes

Solution: Ask Lukas Kahwe Smith



- LAMP software architect
Liip since 2008
- database: MySQL, SQLite,
Oracle, PostGreSQL
- MySQL 5.0 Certified
Developer
- Publications, Speeches and
Contributions in the field of
databases
- www.pooteeweet.org

Lukas K. Smith

Presentation: SQL (un)patterns

Split up into separate tables with as many columns as necessary, create DDL on the fly if necessary

One Table for each content class

content_object_version_myclass1

```
+content_object_id: integer
+version_nr: integer
+locale: string
+name: string
+published: integer
+creator_id: string
+created: timestamp
+additional fields from attributes: mixed
+attrib1_weisnich: integer
+attrib1_blablu: string
+attrib2_sonst
```

content_object_version_myclass2

```
+content_object_id: integer
+version_nr: integer
+locale: string
+name: string
+published: integer
+creator_id: string
+created: timestamp
+additional fields from attributes: mixed
+attrib1_datafield1: integer
+attrib1_datafield2: string
+attrib2_datafield: blob
+attrib3_datafield: timestamp
```

DDL on the fly

How it works

- Adding an attribute: ALTER TABLE ADD COLUMN
- Deleting an attribute: ALTER TABLE DROP COLUMN
- Get the whole object: SELECT *

critique

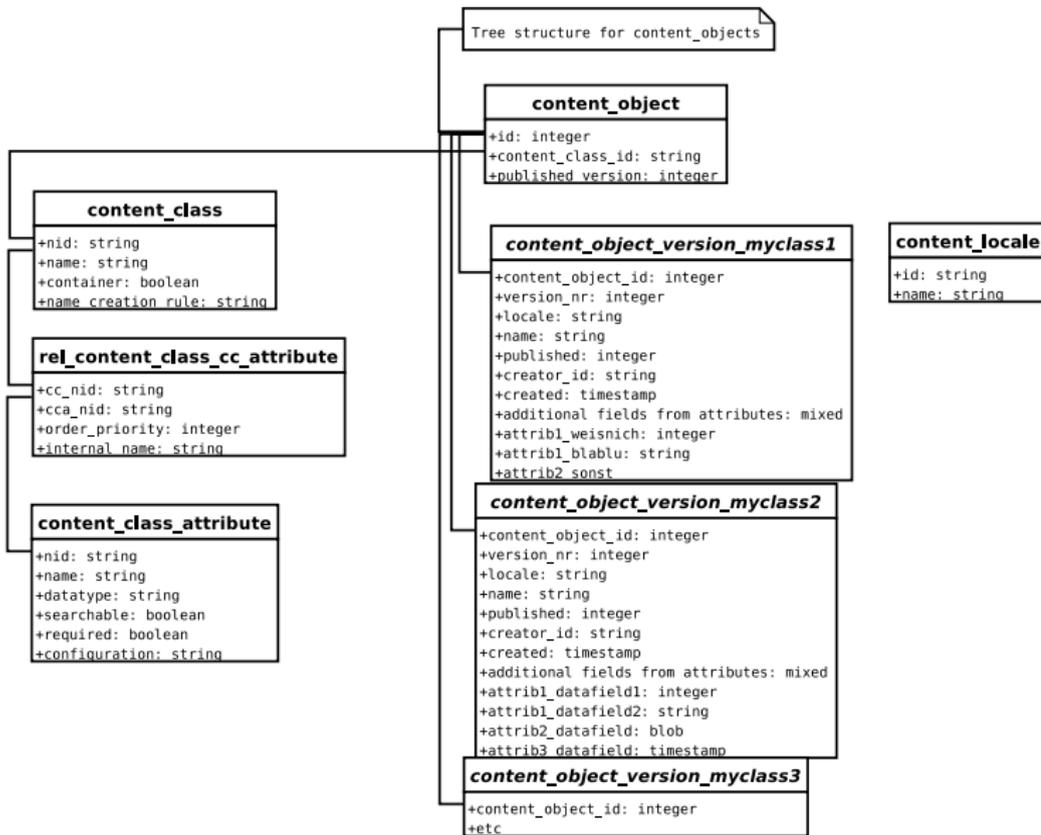
DDL is not transactional save

- MyISAM also does not have transactions and is used in production
- You do not change your content classes weekly but in a staged development process
- The current implementation in eZ Publish is even worse: PHP updates the related attributes until it times out

critique

You must access many tables to get informations on different object types

- Yes, but most of the time you instantiate the objects anyhow and can get the informations in PHP land
- Use views



conclusion

Databases are powerful tools, but todays PHP applications use only a small subset of the available features.

Outline

- 1 Intro
- 2 Static database schemes
- 3 On Flexibility
- 4 Linking objects**
- 5 Keys

Linking is essential

... ERP5 also implements some extra features to enhance programming productivity. Perhaps the most interesting is the concept of relationship managers, which are objects responsible for keeping the relationships between pairs of objects.

1

¹ERP5: Designing for Maximum Adaptability. In: Beautiful Code, S.345

Links in eGroupWare

```
'egw_links' => array(  
  'fd' => array(  
    'link_id' => array('type' => 'auto', 'nullable' => False),  
    'link_app1' => array('type' => 'varchar', 'precision' => '25', 'nullable' => False),  
    'link_id1' => array('type' => 'varchar', 'precision' => '50', 'nullable' => False),  
    'link_app2' => array('type' => 'varchar', 'precision' => '25', 'nullable' => False),  
    'link_id2' => array('type' => 'varchar', 'precision' => '50', 'nullable' => False),  
    'link_remark' => array('type' => 'varchar', 'precision' => '100'),  
    'link_lastmod' => array('type' => 'int', 'precision' => '8', 'nullable' => False),  
    'link_owner' => array('type' => 'int', 'precision' => '4', 'nullable' => False)  )  
)
```

Linking of contents must be a first class priority for the storage layer.

- Put 1:n links inside the table to reduce joins
- Individual n:m link tables for every pair of linked tables?

Outline

- 1 Intro
- 2 Static database schemes
- 3 On Flexibility
- 4 Linking objects
- 5 Keys**

Surrogate Keys vs. Natural Keys

Surrogate Keys

```
'id' int(11) NOT NULL auto_increment
```

Natural Keys

unique identifiers with meaning:

- en, de, ro, fr, ...
- postal codes
- user logins
- email addresses

On surrogate keys

.. They suck.

- surrogate keys give no clue on the object they represent
- most of the times they are another unique index that needs to be checked by the DB
- are not portable from one system to another (dev - production)

Example

Content classes and class attributes in eZPublish have unique string identifiers and surrogate keys.

On natural keys

.. They are lovely.

- gives a hint for debugging and avoids wrong numbers
- often save an id column
- are portable from one system to another (dev - production)

Natural keys may include multiple columns.

on speed

- Don't overoptimize the DB! Do caching!
- Use binary instead of text: No colation
- int takes four bytes and all of them need to be compared, three letters can already be enough to identify a natural key!

links

- <http://archives.postgresql.org/pgsql-hackers/2006-01/msg00414.php>
- <http://pooteweet.org/blog/1015>
- http://en.wikipedia.org/wiki/Surrogate_key