



Web Search By The People For The People

a decentralised search engine:
the missing link
between free content and free users
to build a free web

- **Vision:**

YaCy as the missing link between free content and the user; for privacy, independence, against censoring and for better search results

- **Demonstration:**

what you can do in just five minutes:

installation, crawling, searching, monitoring, scheduling

- **Technology:**

Search engine technology, crawling the web, understanding documents, ranking, peer-to-peer architecture of YaCy and privacy protection

- **Development:**

APIs that you can use with existing tools and easily with some coding.

There is a **missing link** in the web
between **free content** and the **user**

because

free content needs a **decentralised free search** technology



Free Software
Data under Creative Commons License
Open Access Repositories

as it is today:
PROPRIETARY & CENTRALISED:
it **traces you** & data can be **censored, blocked, removed, spammed**

User needs proprietary & centralised software to discover free content

is this what we want?

The World Wide Web should be a **many-to-many media**:
a receiver can be a sender and vice versa



In a **free world wide web** **the users** must run **search engines**

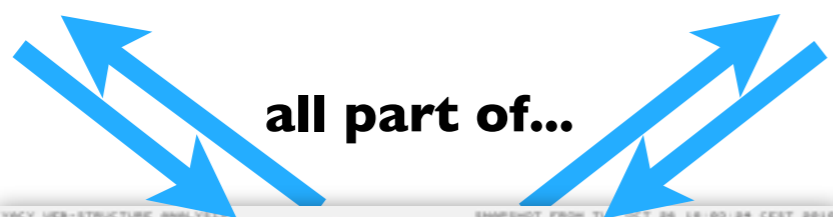
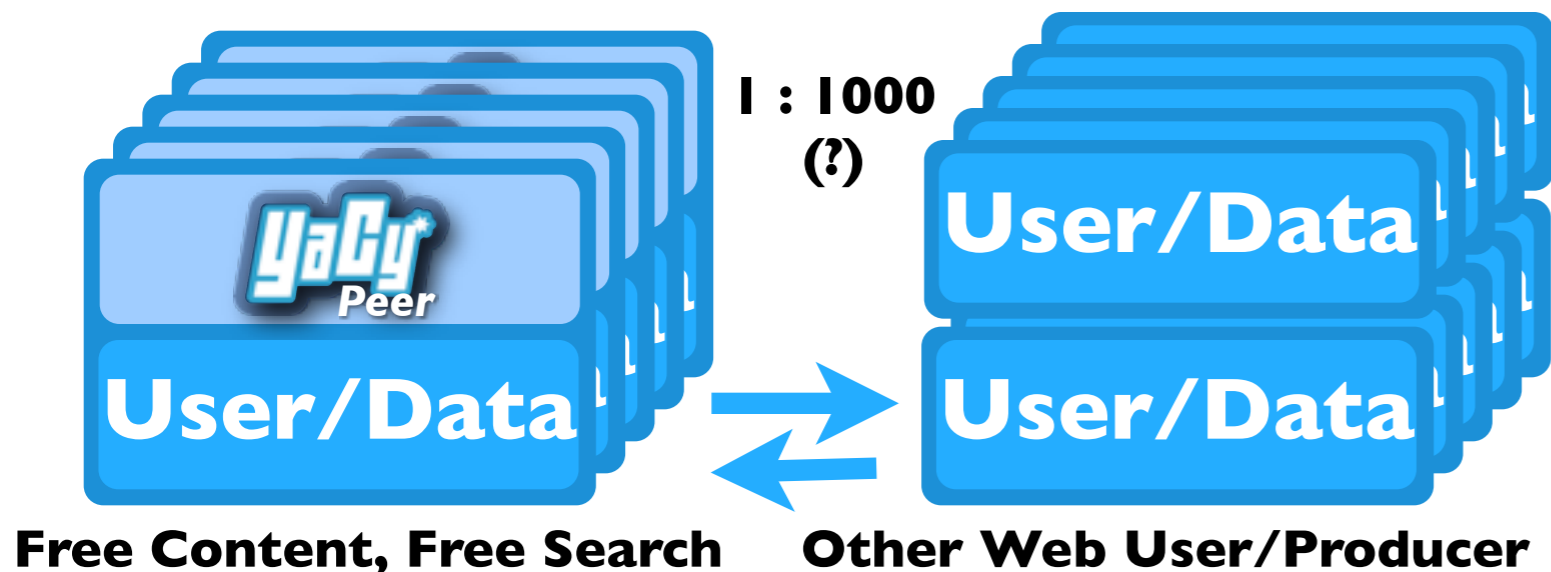


the missing link

between free software, free content and the user
is a network of **YaCy Peers** or something similar



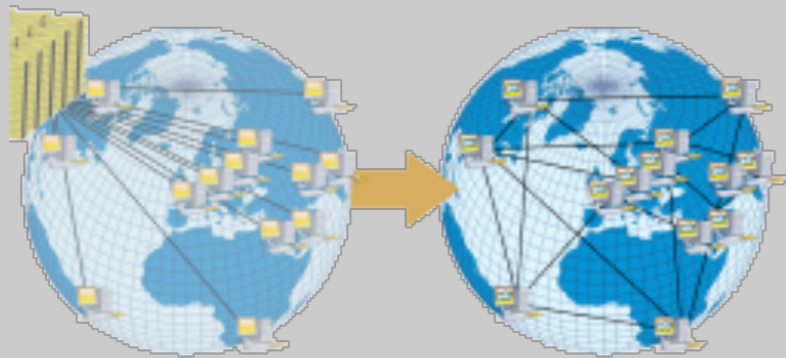
A Free World Wide Web



The Web: Sender & Receiver

Benefits:

- no global censoring*
because its decentralised
- you cannot be traced*
you run the search portal
- same rights for everyone*
everyone can contribute equally
- this is the wiki principle*
for search engines
- get news quickly*
because people care
- better search results*
choose a ranking for your
personal content and relevance



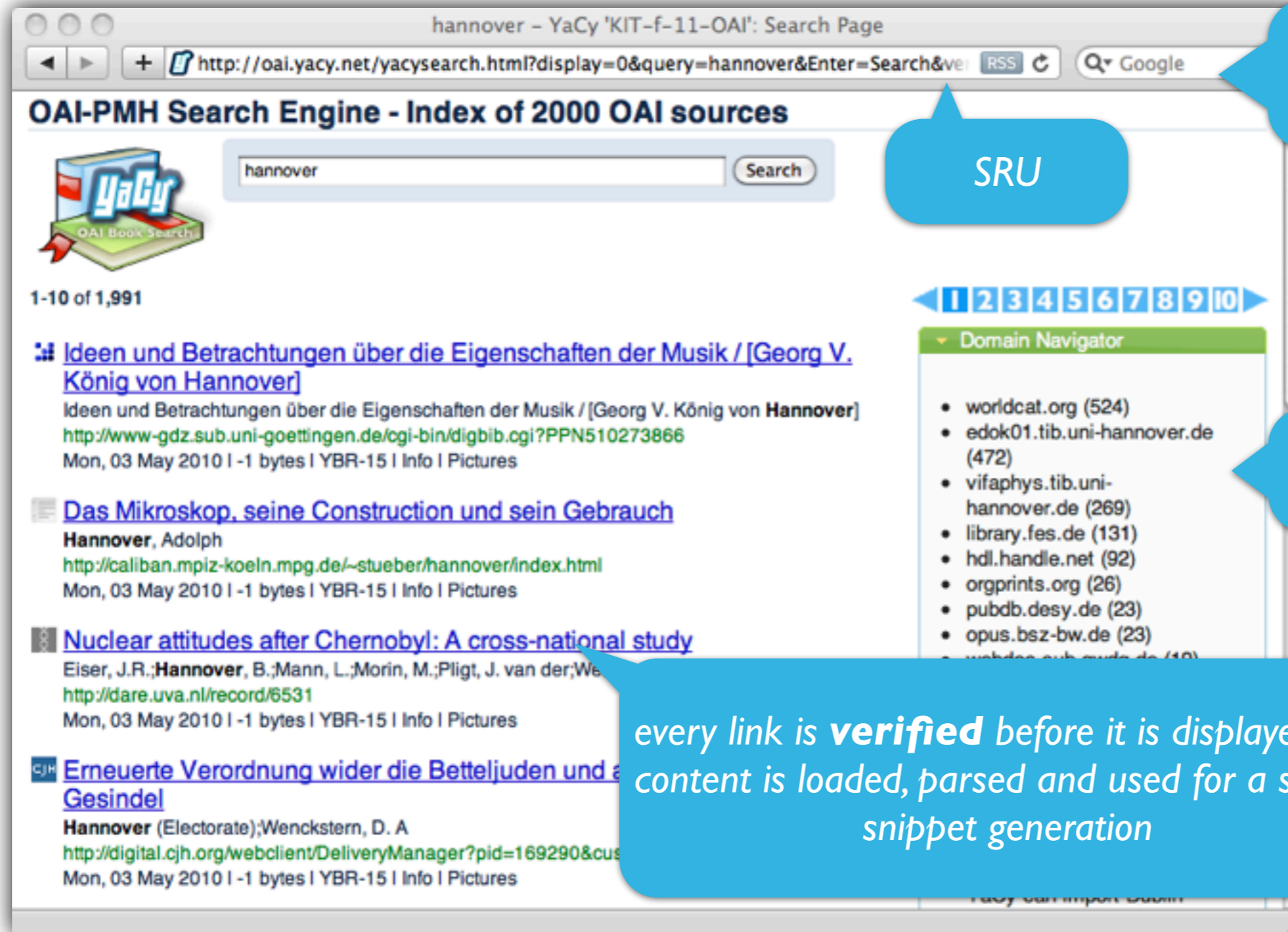
- **Decentralised Peer-to-Peer Web Search**
free search for everyone



- **Internet Search Portal** for a project combining wikis, blogs, forums and portal pages
- **Alert-Service for News using RSS**
create a News-Feed using recent search results for a specific topic



- **Intranet Search Appliance**
search in local web servers and file shares



API for search results is RSS (Opensearch) and JSON

SRU

Facets: Domains, Authors

every link is **verified** before it is displayed: the content is loaded, parsed and used for a search snippet generation

Standards

APIs Opensearch (search results with RSS), JSON, AJAX tools
 Tools search widget, ready-to-use code snippets to embed search everywhere



linuxtag.org

Search Results for linuxtag.org

- Freie Software Presseagentur - LinuxTag
- LinuxTag 2009 - The Lin
- LinuxTag 2009 - Der Lin

linux-club.de

geoclub.de

Geocaching Suchportal

karlsruhe

1-10 aus 4.505

Location - click on map to enlarge

Karlsruhe (lat=49.094, lon=8.5635)

Karlsruhe (Baden) (lat=49.0047, lon=8.3858)

Schnuefflers Geocaching Blog - Karlsruhe

Gefunden: Alien in Karlsruhe

wegeundpunkte ... Rheinhafen Karlsruhe

Karlsruhe - Von Hamburg bis Hawaii

GEOCACHING GPS - Karlsruhe

fsfe.org

Free Software Foundation Europe Search Portal

lizenz

1-10 of 34

FSF Europe - WIPO-Beobachtung - Referenzblatt zu den Grundlagen Freier Software

Die FTF und gpl-violations.org veröffentlichen eine Anleitung zum Melden und Beheben von Lizenz-Verletzungen

FSFE - Freedom Task Force (FTF)

FSFE - Freedom Task Force (FTF)

Search results provided by YaCy

YaCy Portal Search

YACY DISTRIBUTED WEB SEARCH

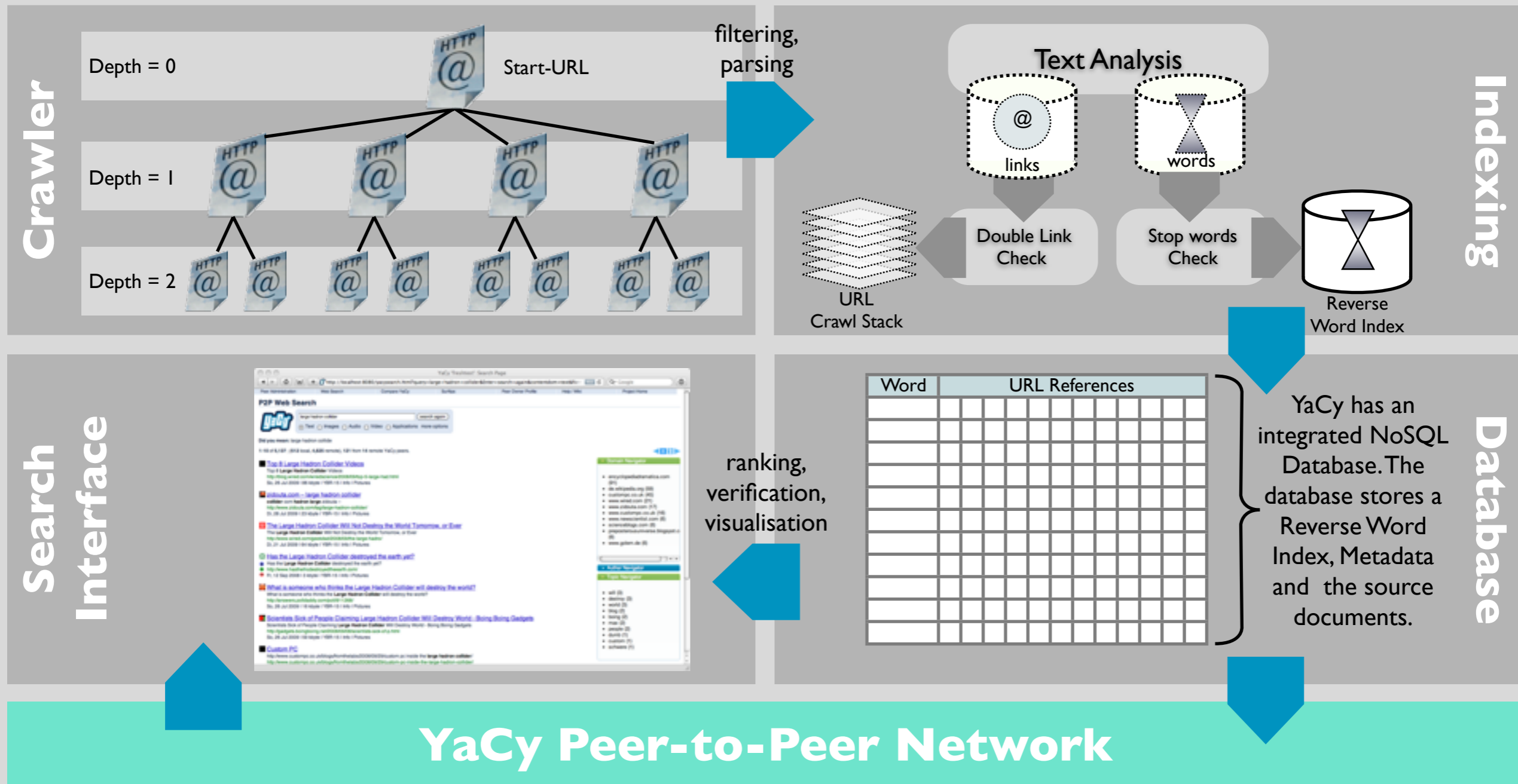
DOMAINS

- www.linux-club.de (2915)
- www.learninglinux.de (579)
- www.linuxtag.org (346)
- www.netsecond.net (82)
- www.kefk.net (73)
- www.fini-online.com (25)
- www.linuxmintusers.de (5)
- www.ubuntu-freunde.de (4)
- yacy.net (2)
- www.golem.de (2)

AUTHORS

TOPICS

Retrieval, Indexing, Storage and Search Components



Crawler Administration

Site Crawl Start

Site Start URL

Link-List of URL
 Sitemap URL

Scheduler run this crawl once
 scheduled, look every days for new documents automatically.

Path load all files in domain
 load only files in a sub-path of given url

Limitation not more than documents

Dynamic URLs allow query-strings (urls with a '?' in the path)

Start

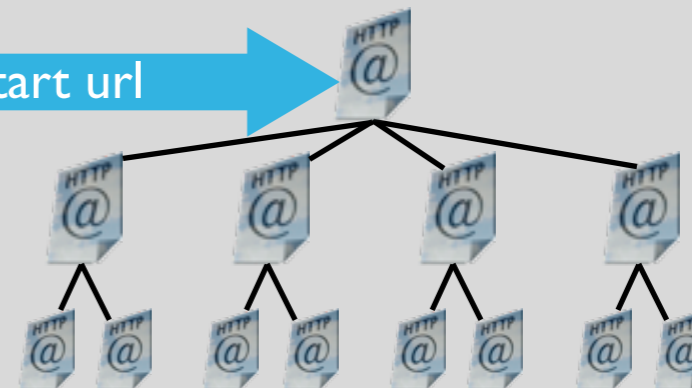
a simple 'Site Crawl'; there is also a detailed crawl start for 'wide' crawls

all crawl starts are placed into a scheduler

Type	Comment	Call Count	Last Exec Date	Next Exec Date	Scheduler
<input type="checkbox"/> crawler	crawl start for http://blogs.fsfe.org/	50	Oct 31, 2010 1:50:18 PM	Nov 1, 2010 1:49:00 PM	<input type="text" value="1"/> days
<input type="checkbox"/> crawler	crawl start for http://download.fsfe.org/	8	Oct 25, 2010 2:50:57 PM	Nov 1, 2010 1:49:00 PM	<input type="text" value="7"/> days
<input type="checkbox"/> crawler	crawl start for http://fellowship.fsfe.org/	8	Oct 25, 2010 2:50:57 PM	Nov 1, 2010 1:49:00 PM	<input type="text" value="7"/> days
<input type="checkbox"/> crawler	crawl start for http://fsfe.org/	8	Oct 25, 2010 2:50:57 PM	Nov 1, 2010 1:49:00 PM	<input type="text" value="7"/> days
<input type="checkbox"/> crawler	crawl start for http://lists.fsfe.org/	8	Oct 25, 2010 2:50:57 PM	Nov 1, 2010 1:49:00 PM	<input type="text" value="7"/> days
<input type="checkbox"/> crawler	crawl start for http://planet.fsfe.org/	8	Oct 25, 2010 2:50:57 PM	Nov 1, 2010 1:49:00 PM	<input type="text" value="7"/> days
<input type="checkbox"/> crawler	crawl start for http://wiki.fsfe.org/	26	Oct 31, 2010 1:50:18 PM	Nov 2, 2010 1:49:00 PM	<input type="text" value="2"/> days
<input type="checkbox"/> crawler	crawl start for http://www.drm.info/	8	Oct 25, 2010 2:50:57 PM	Nov 1, 2010 1:49:00 PM	<input type="text" value="7"/> days
<input type="checkbox"/> crawler	crawl start for http://www.pdfreaders.org/	8	Oct 25, 2010 2:50:57 PM	Nov 1, 2010 1:49:00 PM	<input type="text" value="7"/> days
<input type="checkbox"/> crawler	crawl start for http://mail.fsfeurope.org/mailman/listinfo	8	Oct 25, 2010 2:50:57 PM	Nov 1, 2010 1:49:00 PM	<input type="text" value="7"/> days

Document Tree

the start url



follow all links until:

- no more documents in domain
- crawl depth is reached
- maximum number of docs reached

use

Target Host Balancing

(for several crawl starts or 'wide' crawls)

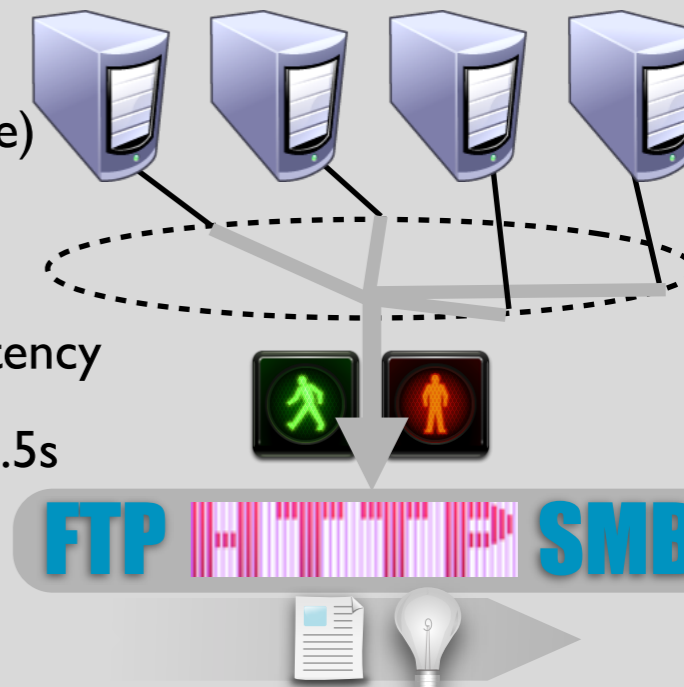
target hosts (domain name)

round-robin access

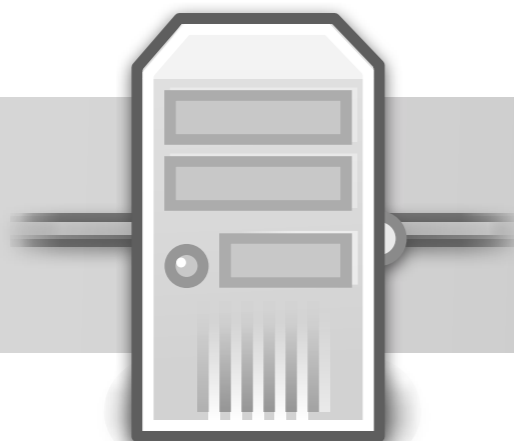
robots.txt, latency and minimum access time 0.5s

loader

parser



A search engine should support **people** in the search for documents in **unstructured** formats: this needs a kind of **'understanding'** of content



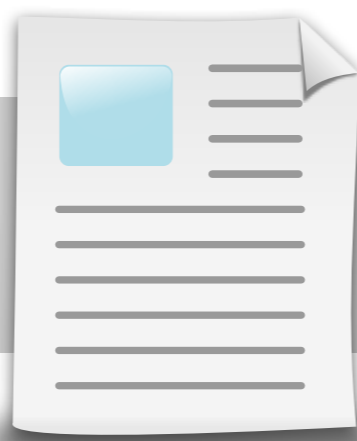
Connection

load and crawl from:

HTTP, HTTPS, FTP, filesystem,
SMB-shares

Import from:

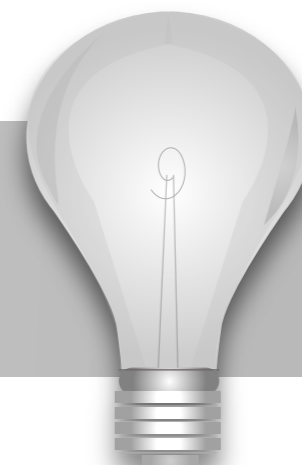
Dublin Core / XML files,
OAI-PMH, wikimedia dumps,
SQL databases



Parsing

read document formats:

HTML, XHTML, RSS, RDF,
XHTML+RDFa, FOAF, vCard,
Flash, PDF, PS, Word, Excel,
Visio, Powerpoint,
OpenOffice, RTF, csv, gzip, zip,
tar, rar, bzip2, 7zip, images
(EXIF), torrent files



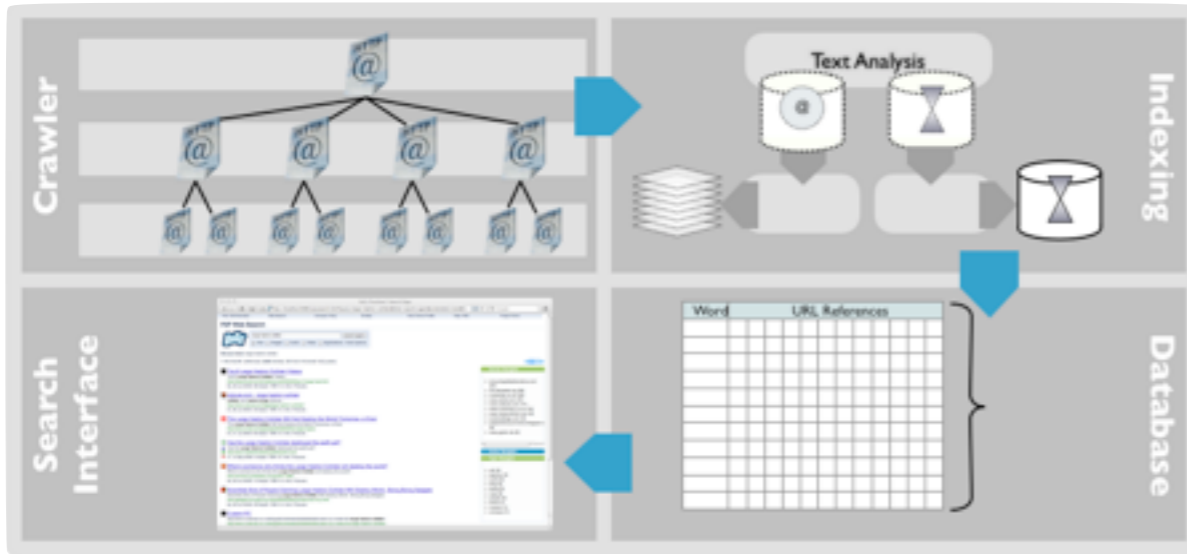
Interpretation

find metadata (headline,
author, date, locations)

find links of different kind
(text, images, movies etc.)

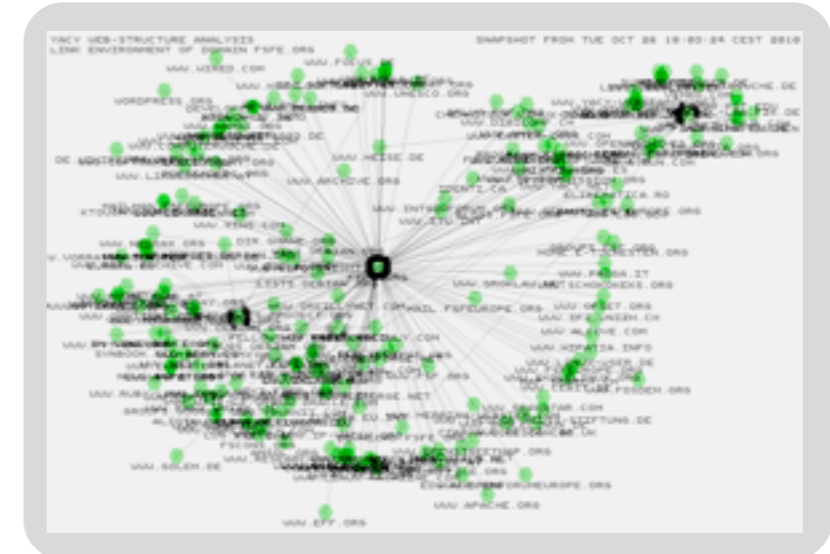
store statistical data for
search suggestions

Search Engine



retrieval, indexing, storage and search components

Data Visualisation



index creation process, system load, link structure, p2p net configuration

Scheduler and Steering



automatic scheduled re-indexing and back-up of search appliance set-up

Database Administration

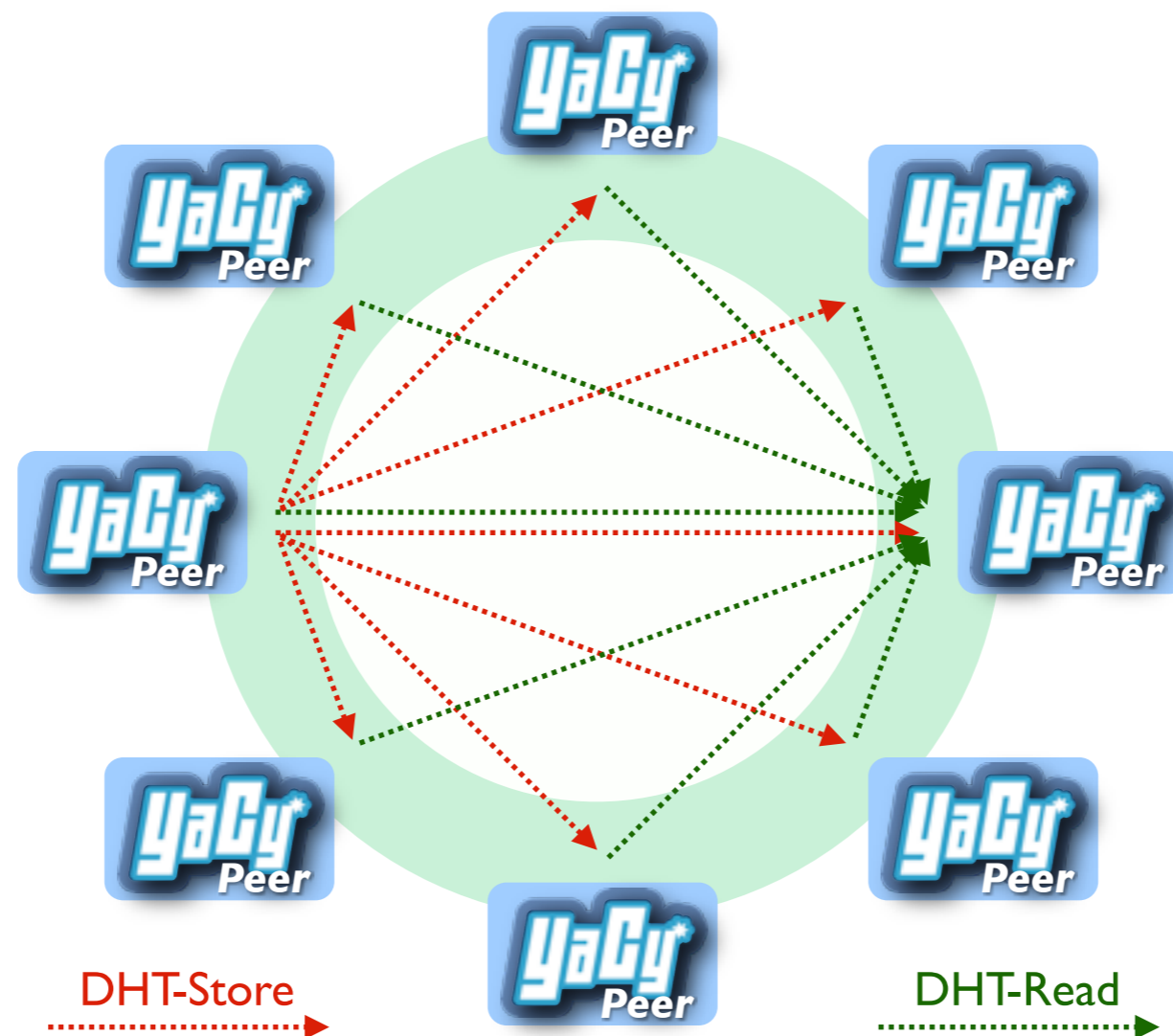
Table Editor: showing table 'rss'

PK	title	recording_date	url
<input type="checkbox"/> -0d90MM7P7UE	Lugares RSS	20100827075402810	http://www.bwnpatagonia.com.ar/feed/lugares.xml
<input type="checkbox"/> -lhCeMM7P7UE	Culturales RSS	20100827075402810	http://www.bwnpatagonia.com.ar/feed/culturales.xml
<input type="checkbox"/> 0LCJL7K-wcia	everyfoodfits.com Posts RSS feed	20101008121101879	http://everyfoodfits.com/feed/
<input type="checkbox"/> 1zyuvMM7P7UE	Ciencia y Tec RSS	20100827075402809	http://www.bwnpatagonia.com.ar/feed/ciencia%20y%20tec.xml
<input type="checkbox"/> 23p96nWc9erZ	Recipes from the Restaurant Refugee RSS Feed	20100910183339465	http://rrecipes.wordpress.com/feed/
<input type="checkbox"/> 2fzNTMM7P7UE	Causas RSS	20100827075402803	http://www.bwnpatagonia.com.ar/feed/causas.xml

Buttons: Edit Selected Row, Add a new Row, Delete Selected Rows, Delete Table

crawl queues, robots.txt, rss feeds, scheduler data, p2p connections, network messages

The YaCy Network: a distributed hash table



This peer (as an example) fetches some Web pages and distributes index fragments to other peers.

A peer which searches information can access directly peers holding the corresponding index

YaCy peers store index fragments according to a 'folded' ordering on word-hashes and url-hashes in a distributed hash table (DHT). The index is distributed redundantly to save the index when some peers are not available. The redundancy also helps to increase search performance.

The Architecture of YaCy Ensures Privacy

Nobody can see what you added to the global search index

- YaCy does not store words in clear text but only as word-hashes
- your search index is mixed with indexes from other peers during DHT transfer
- the origin of DHT transfers is not stored into the search index

Nobody can see what you search

- if you use your own YaCy application, you are the only one who can track what you do
- a tracking against mis-use (DoS etc.) is build-in, but
- remote search tracking cannot see the remote users search words, only word hashes.

How to integrate a YaCy Search Portal:
Just copy-paste the code snippet to your web page source code.

Code Snippet Example #1: a search window in an iframe

```
<iframe name="target2"
  src="http://141.52.175.43:8080/yacyssearch.html?
display=2&resource=local"
  width="100%" height="180"
  frameborder="0" scrolling="auto" id="target2"
</iframe>
```

Code Snippet Example #2: a search box (points to new page)

```
<form method="get" accept-charset="UTF-8"
  action="http://141.52.175.43:8080/yacyssearch.html">
  <div>
    <div>MySearch</div>
    <input type="text" name="query" value="" maxlength="80" />
    <input type="hidden" name="verify" value="true" />
    <input type="hidden" name="maximumRecords" value="10" />
    <input type="hidden" name="meanCount" value="5" />
    <input type="hidden" name="resource" value="local" />
    <input type="hidden" name="urlmaskfilter" value=".*" />
    <input type="hidden" name="prefermaskfilter" value="" />
    <input type="hidden" name="display" value="2" />
    <input type="hidden" name="nav" value="all" />
    <input type="submit" name="Enter" value="Search" />
  </div>
</form>
```

Code Snippet #2 looks like:

MySearch

The YaCy administration interface offers more code snippets. An example from

`/ConfigSearchBox.html`

looks like:

MySearch

your YaCy peer provides help pages with code snippets for an easy integration!



```
> curl http://localhost:8080/yacyssearch.rss?query=foaf&maximumRecords=10
<?xml version="1.0" encoding="UTF-8"?>
<?xml-stylesheet type='text/xsl' href='/yacyssearch.xsl' version='1.0'?>
<rss version="2.0"
  xmlns:yacy="http://www.yacy.net/"
  xmlns:opensearch="http://a9.com/-/spec/opensearch/1.1/"
<!-- very short example -->
<item>
  <title>Friend of a Friend (FOAF) project</title>
  <link>http://www.foaf-project.org/</link>
  <pubDate>Fri, 23 May 2008 02:00:00 +0200</pubDate>
</item>
<item>
  <title>FOAF - Wikipedia</title>
  <link>http://de.wikipedia.org/wiki/FOAF</link>
  <pubDate>Tue, 08 Jan 2008 01:00:00 +0100</pubDate>
</item>
<item>
  <link>http://microformats.org/wiki/xfn-to-foaf</link>
  <pubDate>Fri, 09 May 2008 02:00:00 +0200</pubDate>
</item>
</rss>
```

How to get Opensearch/JSON Search Results:

- do a normal web search in YaCy
- replace the 'html' extension of the result page URL with 'rss'
- for json, replace the 'html' extension with 'json'

Standards:

The YaCy-internal Dublin Core Metadata Format fits into the RSS format for search result data in Opensearch standard very well.

If wanted, also JSON can be used as export format.

Opensearch Standard: <http://www.opensearch.org>

SRU Standard for Queries: <http://www.loc.gov/standards/sru/specs/search-retrieve.html>

```
<?xml version="1.0" encoding="utf-8"?>
<!-- YaCy surrogate using dublin core notion -->
<surrogates
  xmlns:dc="http://purl.org/dc/elements/1.1/">

  <record>
    <dc:title><![CDATA[Alan Smithee]]></dc:title>
    <dc:identifier>http://de.wikipedia.org/wiki/Alan_Smithee</dc:identifier>
    <dc:description>
      <![CDATA[''Alan Smithee'' ist ein Anagramm von „The Alias Men“.]>
    </dc:description>
    <dc:language>de</dc:language>
    <dc:date>2009-04-14T00:00:00Z</dc:date>
    <!-- date is in ISO 8601 -->
  </record>
</surrogates>
```

Standards:

YaCy can import standard Dublin Core Metadata XML files as input for indexing

How to import Dublin Core Files:

just place the xml files into a hand-over directory at DATA/SURROGATES/in/

The Dublin Core XML File Standard:

<http://dublincore.org/documents/dc-xml-guidelines/>

- **Download from <http://yacy.net>**



YaCy

YaCy for Windows



YaCy.app

YaCy for Mac

`yacy_0.97svn7217_all.deb`

YaCy for Debian

`yacy_v0.97_20101004_7217.tar.gz`

YaCy for Linux / generic (tar.gz)

- **Just Extract the Package, then Start the Start-Script**

There are simple installers for Windows, Mac and a debian release, but it is easy to just install the generic release because it contains everything that is needed.

- **Administration using the Web Interface**

YaCy is a Web Application. The administration can be done completely using the built-in web interface with your web browser. Just open `http://localhost:8080`

The main configuration is done when you select your use case (Distributed P2P Web Search, Portal Search, Intranet Search) after just two clicks.

- **Support**

We have a web forum: `http://forum.yacy.de`

Some information can be found at the wiki: `http://wiki.yacy.de`

...or contact me: `mc@yacy.net`



Where is a (demo) Search Portal?

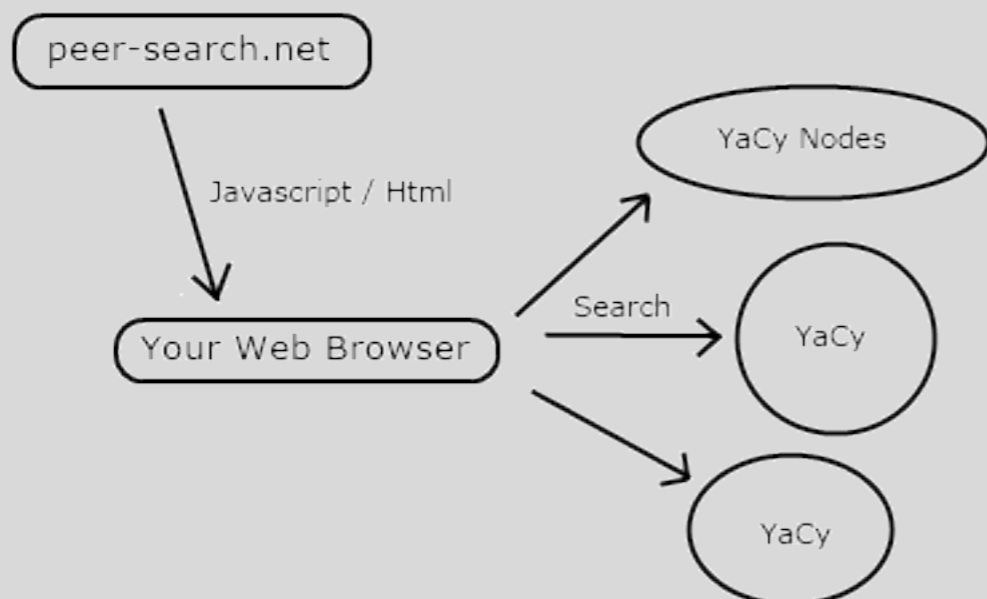
There is no one-for-all demo portal for YaCy!

YaCy is about decentralised search and offering a central point for everyone would ruin the idea!

Distributed Search in Your Browser:

<http://peer-search.net>

- JavaScript Code is loaded into your browser
- your browser loads a list of YaCy peers
- when you search, your browser contacts some of the YaCy peers and combines the search results from these peers; like a meta-search.



Peer Roulette, search on a random peer:

<http://www.yacyweb.de/peers.htm>

- yacyweb generates a list of YaCy peers
- when you click on a link you get the web interface of the peer directly
- when you search on that peer the content may be restricted to the rules of the peer owner

Active YaCy-Peers

Phoenix	pipapo	KIT01-01-FESEARCHING
dlc-TEOx	osuchardotnet	senior
suma-ev2	Lafkor	Eechum
KIT01-03-checker	mgs-athvox	KIT01-02-2008103-FSEE
whitecloud	free-search	yacystats.de-01
rrznhd	w-sci	ImageBroWurst
suma-ev1	Agadius	0x001e
sixcooler	debian-suche2	dlc-heday
rrzn1	Fastbull21	sender
konstantin@homedns.org	connect0a	dortmund-initiativ
nyduzifkpmoxone	KIT01-03-FELINKUREDA	sixcooler2

The best demo: run your own peer!

please help to create
a free web

- *tweet or blog about YaCy (or re-tweet yacy_search)*
- *run and use your own peer*
- *use YaCy to create your own search portal*
- *this is search engine research: tell us your ideas!*
- *this is free software: please submit your code!*

Thank You!