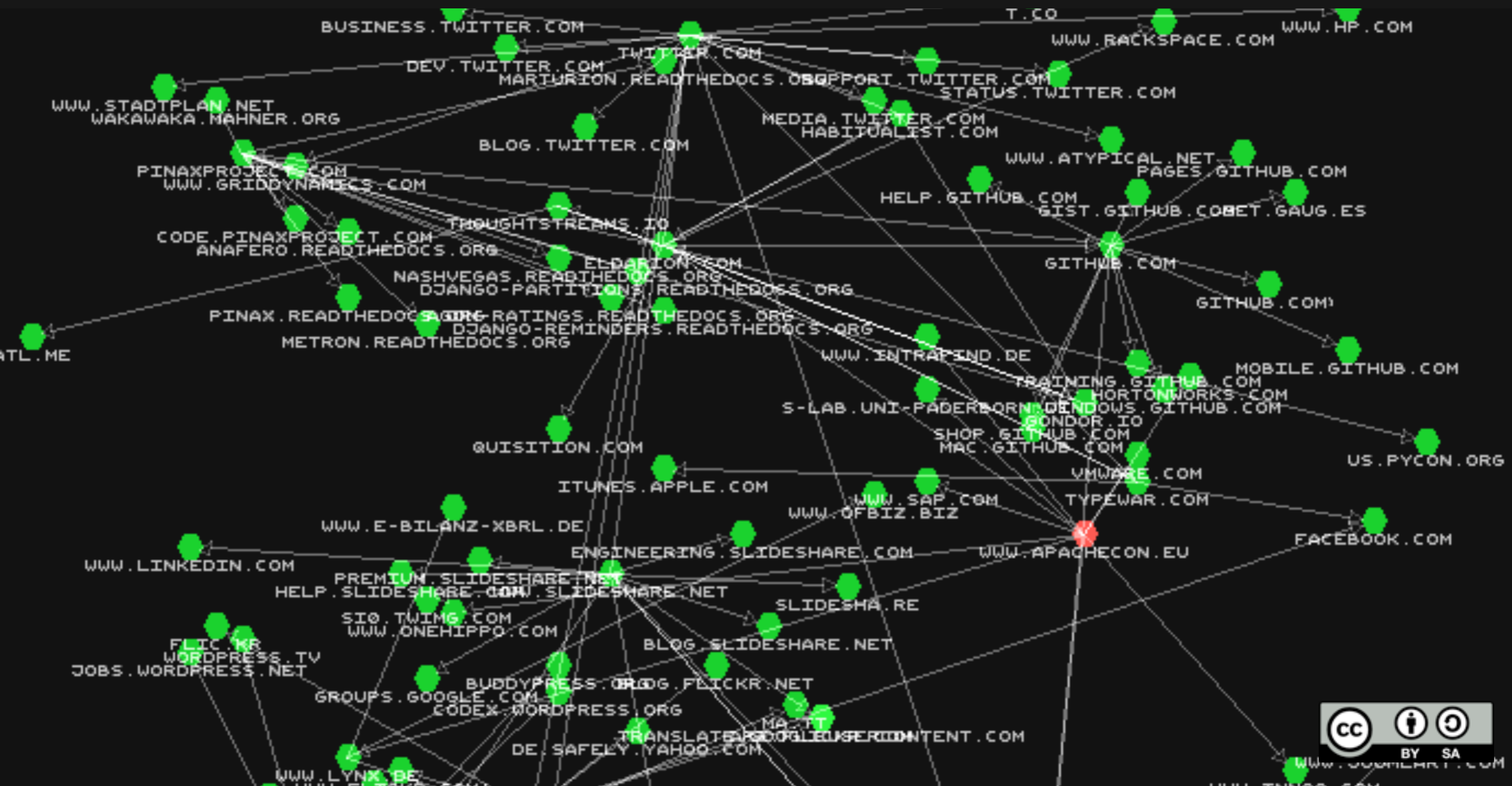




Information Retrieval Web Crawling mit YaCy

Michael Christen mc@yacy.net



URLs pro Woche, deren Löschung aus der Suche beantragt wurde



Im letzten Monat für die Suche eingegangene Ersuchen um Löschung von urheberrechtsverletzenden Inhalten

35.449.471	URLs, deren Löschung beantragt wurde
56.317	Angegebene Domains
4.897	Urheberrechtsinhaber
2.126	Ersuchende Unternehmen

<http://www.google.com/transparencyreport/removals/copyright/?hl=de>

Stand: Januar 2015

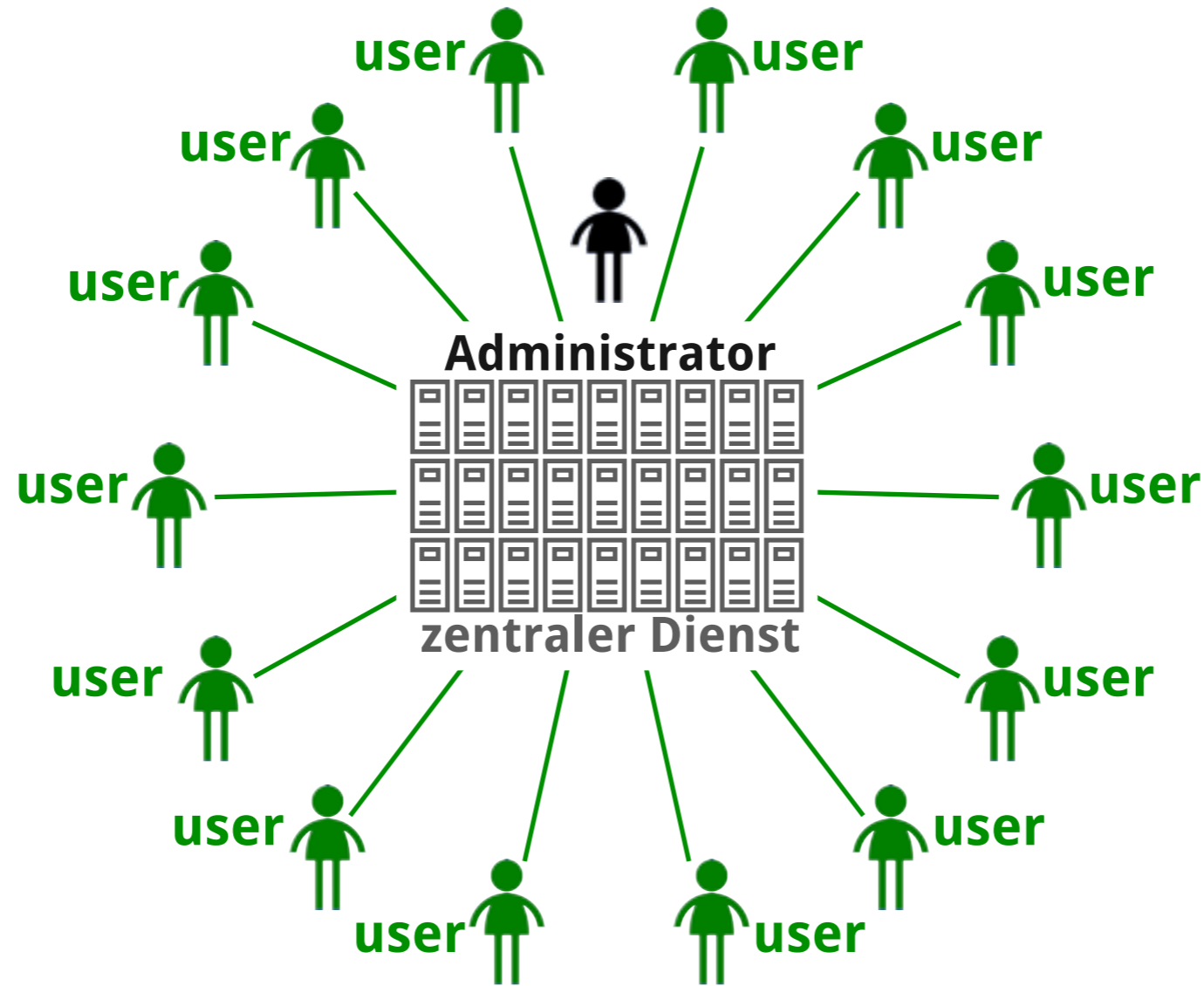
Ersuchen zur Löschung von Suchergebnissen gemäß europäischem Datenschutzrecht

www.facebook.com URLs entfernt: 5126	profileengine.com URLs entfernt: 4901	groups.google.com URLs entfernt: 3446	badoo.com URLs entfernt: 3168	www.youtube.com URLs entfernt: 3082
www.wherevent.com URLs entfernt: 2133	www.yasni.de URLs entfernt: 2128	www.192.com URLs entfernt: 2071	plus.google.com URLs entfernt: 1928	www.yasni.fr URLs entfernt: 1786

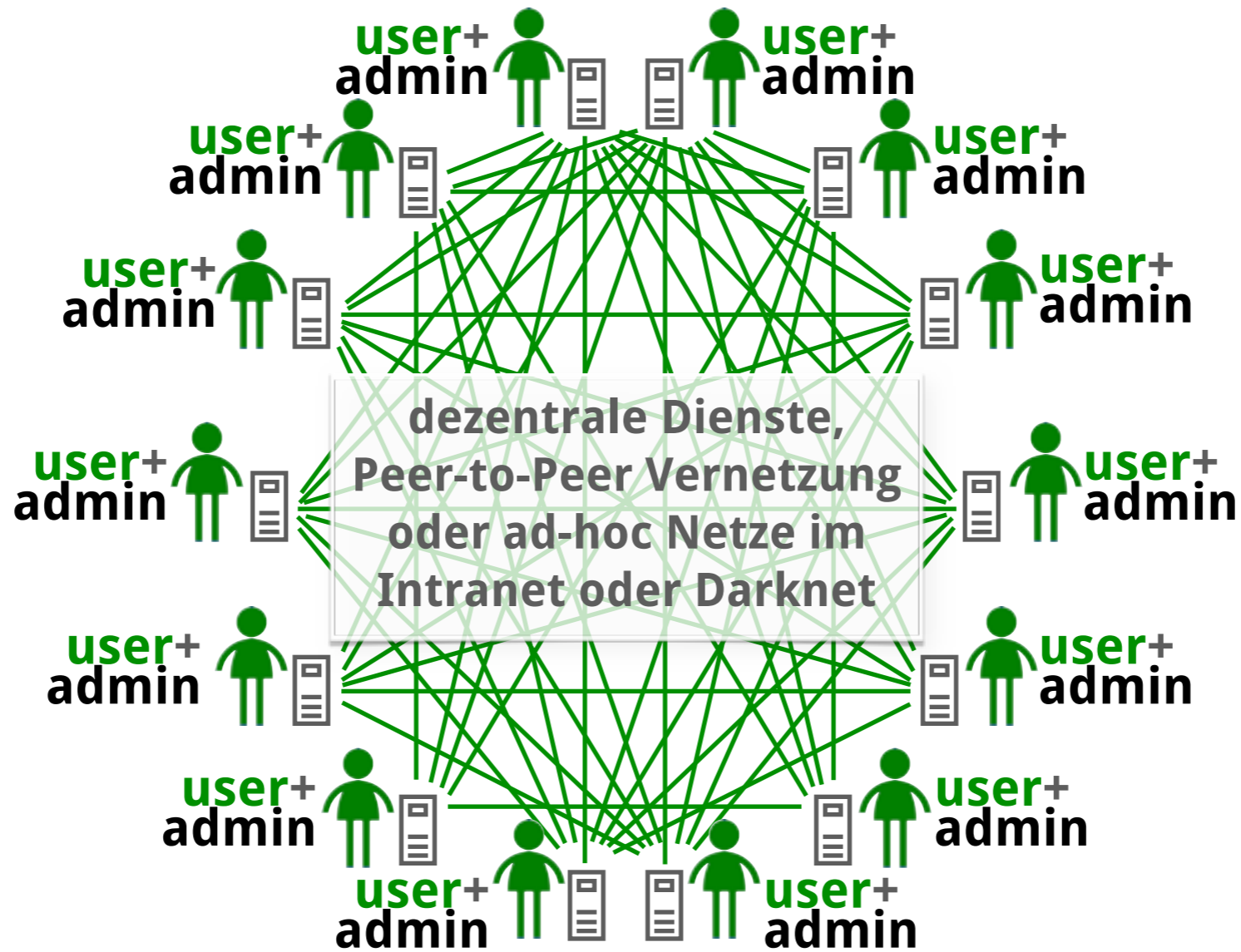
<http://www.google.com/transparencyreport/removals/europeprivacy/?hl=de>

Stand: Januar 2015

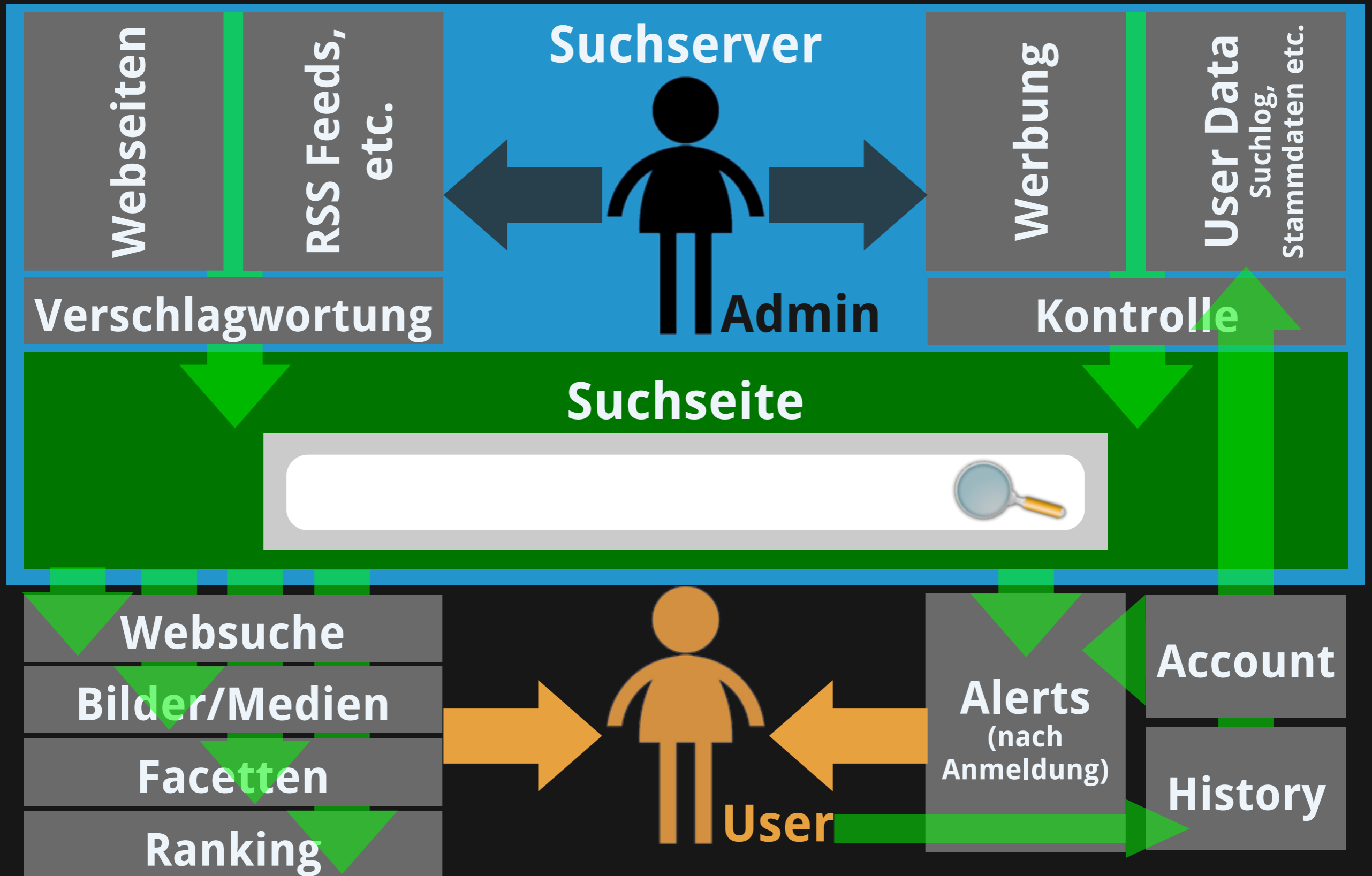
Konzept: Dezentralisierung



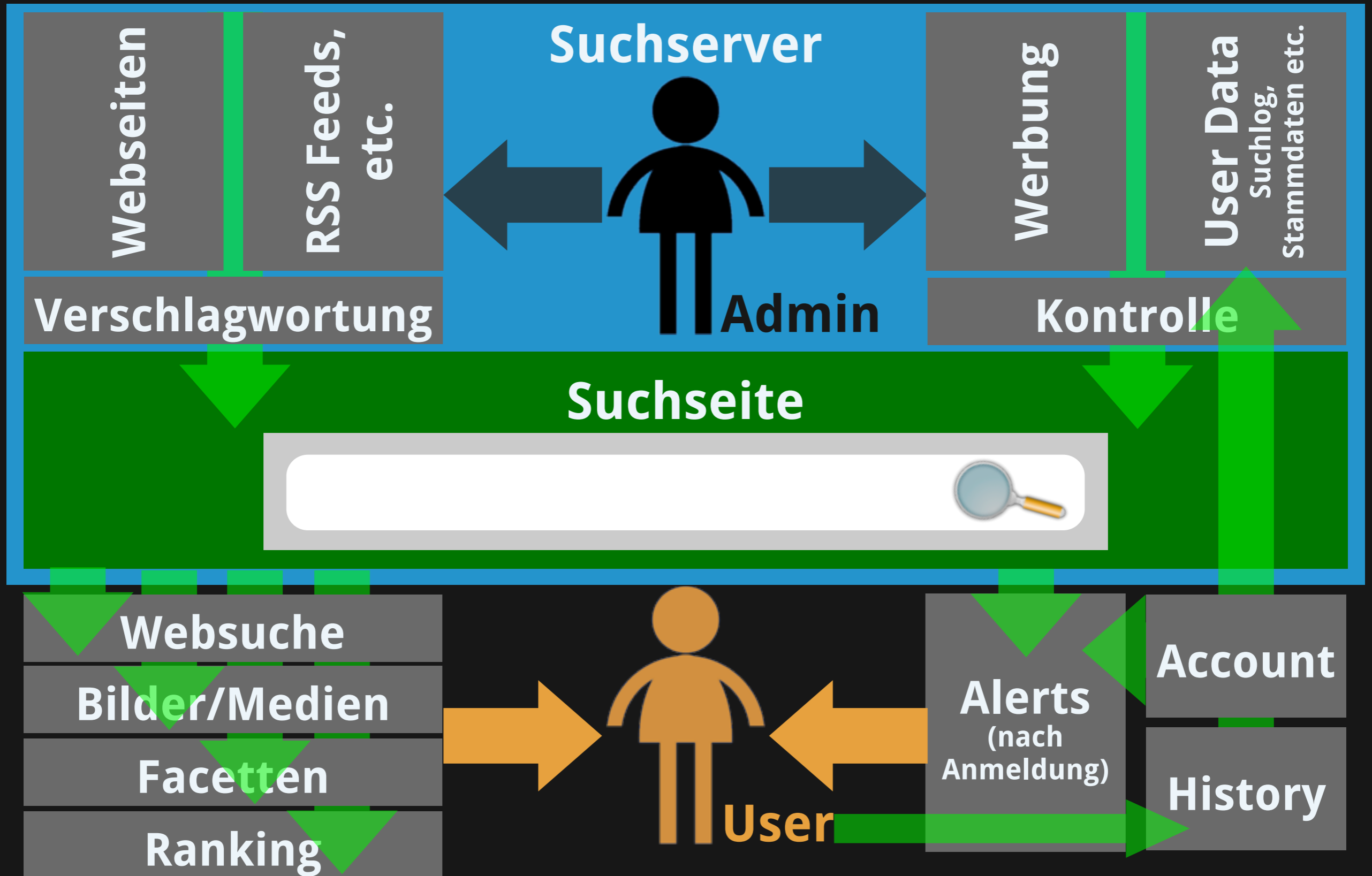
Konzept: Dezentralisierung



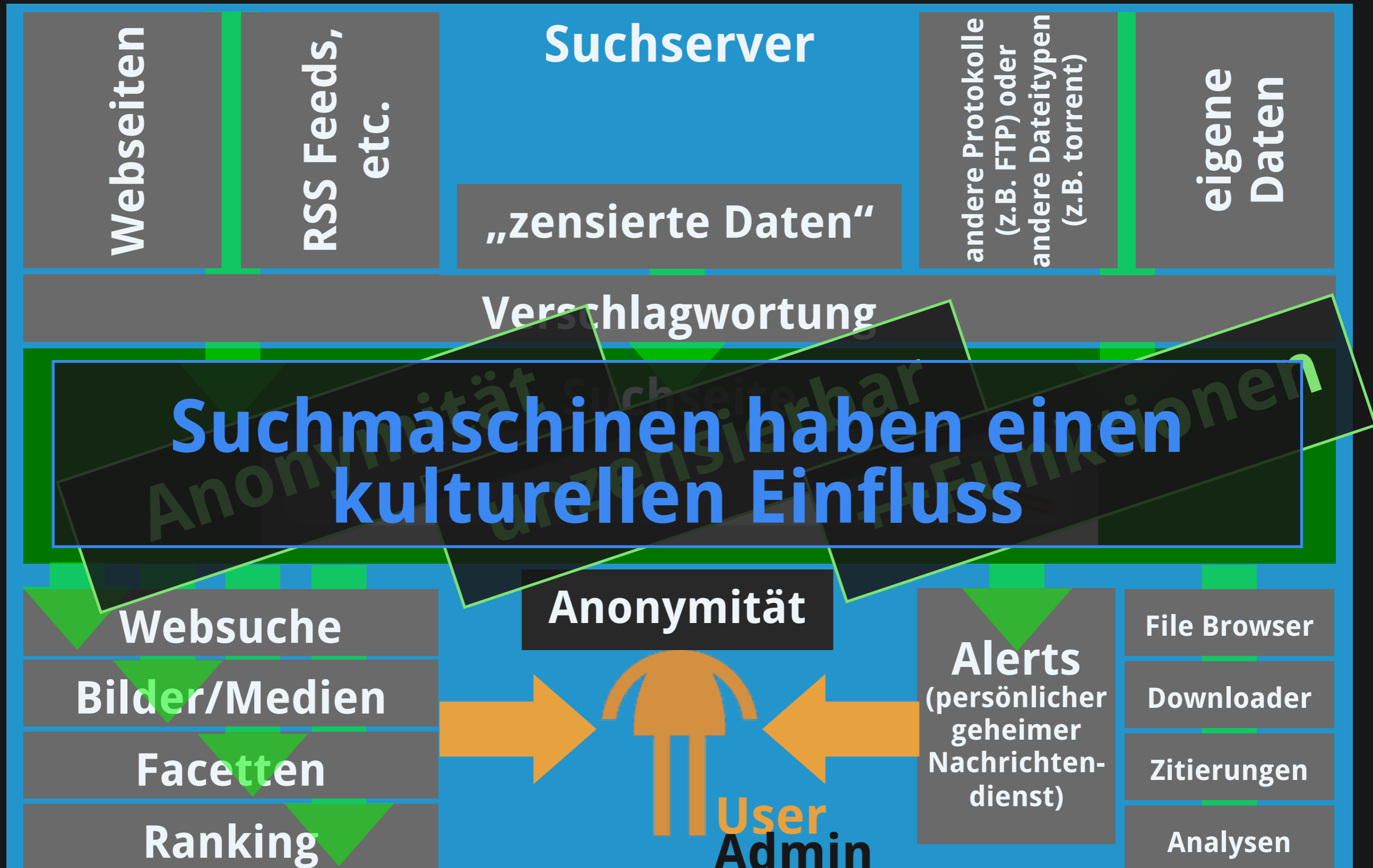
Wissen wie Suchmaschinenportale funktionieren

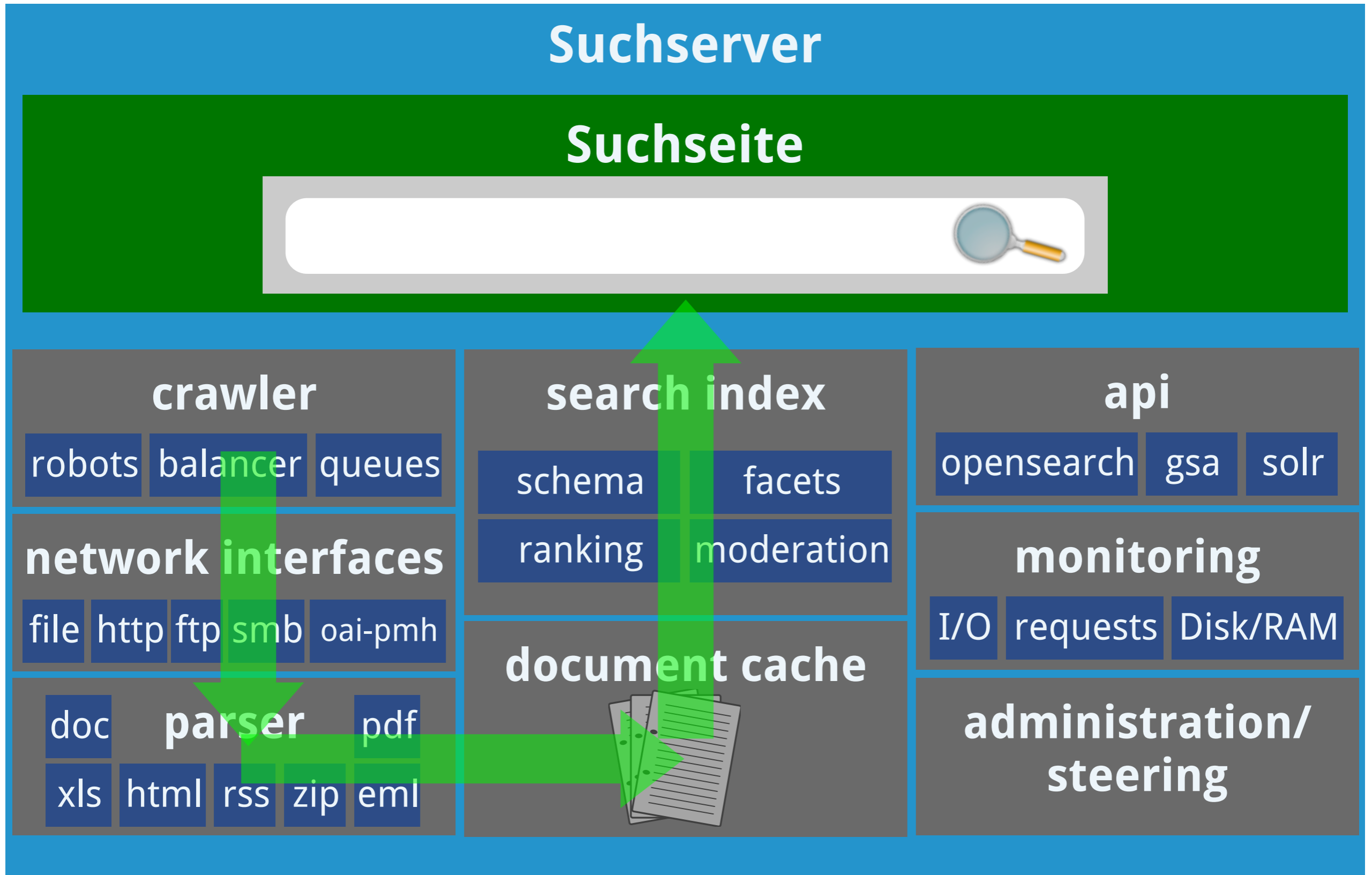


Wunsch überflüssiges entfernen



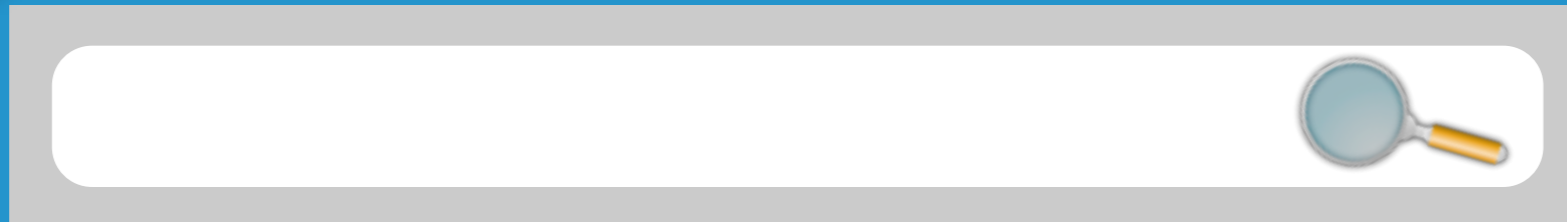
Selbstgemachte Suchmaschine







Suchserver

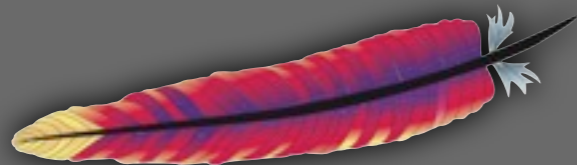


crawler

search index

api

network interfaces

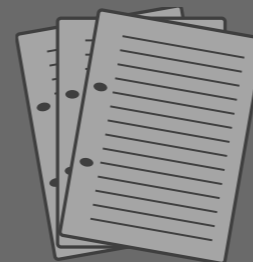


monitoring

parser



document cache



**administration/
steering**

Funktionskomponenten einer Suchmaschine



Komponenten der Suchmaschine

crawler



Komponenten der Suchmaschine

crawler



Crawl Job

A Crawl Job consist of one or more start point, crawl limitations and document freshness rules.

Start Point

One Start URL or a list of URLs:
(must start with http:// https:// ftp:// smb:// file://)

From Link-List of URL
From Sitemap

From File (enter a path within your local file system)

Crawler Filter

Crawling Depth

also all linked non-parsable documents

Unlimited crawl depth for URLs matching with

Maximum Pages per Domain

Use: Page-Count:

misc. Constraints

Accept URLs with query-part ("?"): Obey html-robots-noindex:

Load Filter on URLs

+ must-match

Restrict to start domain(s)

Restrict to sub-path(s)

Use filter

(MUST not be empty)

- must-not-match

Clean-Up before Crawl Start

No Deletion

Do not delete any document before the crawl is started.

Delete sub-path

For each host in the start url list, delete all documents (in the given subpath) from that host.

Delete only old

Treat documents that are loaded days ago as stale and delete them before the crawl is started.

Double-Check Rules

No Doubles

Never load any page that is already known. Only the start-url may be loaded again.

Re-load

Treat documents that are loaded days ago as stale and load

Funktionskomponenten einer Suchmaschine

Komponenten der Suchmaschine

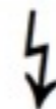
crawler

parser



Content Parser Settings

enable/disable	Extension
<input checked="" type="checkbox"/>	Microsoft Powerpoint Parser
<input checked="" type="checkbox"/>	pps
<input checked="" type="checkbox"/>	ppt
<input checked="" type="checkbox"/>	GNU Zip Compressed Archive Parser
<input checked="" type="checkbox"/>	gz
<input checked="" type="checkbox"/>	tgz
<input checked="" type="checkbox"/>	Adobe Flash Parser
<input checked="" type="checkbox"/>	swf
<input checked="" type="checkbox"/>	vCard Parser
<input checked="" type="checkbox"/>	vcf
<input type="checkbox"/>	Audio File Meta-Tag Parser
<input type="checkbox"/>	m4p
<input type="checkbox"/>	m4a
<input type="checkbox"/>	oga
<input type="checkbox"/>	flac
<input type="checkbox"/>	ogg
<input type="checkbox"/>	mp3
<input type="checkbox"/>	wma
<input type="checkbox"/>	Comma Separated Value Parser
<input type="checkbox"/>	csv
<input type="checkbox"/>	Microsoft Visio Parser
<input checked="" type="checkbox"/>	vdx
<input checked="" type="checkbox"/>	vtx
<input checked="" type="checkbox"/>	vss
<input checked="" type="checkbox"/>	vsd
<input checked="" type="checkbox"/>	vst
<input checked="" type="checkbox"/>	Generic Image Parser
<input checked="" type="checkbox"/>	bmp
<input checked="" type="checkbox"/>	jpg
<input checked="" type="checkbox"/>	jpeg
<input checked="" type="checkbox"/>	png
<input checked="" type="checkbox"/>	jpe
<input checked="" type="checkbox"/>	gif
<input checked="" type="checkbox"/>	FreeMind Parser
<input checked="" type="checkbox"/>	mm
<input checked="" type="checkbox"/>	PostScript Document Parser
<input checked="" type="checkbox"/>	ps
<input checked="" type="checkbox"/>	Commodore 64 SID Audio File Parser
<input checked="" type="checkbox"/>	sid



Open Office XML Document Parser

<input checked="" type="checkbox"/>	xltx
<input checked="" type="checkbox"/>	xlsx
<input checked="" type="checkbox"/>	ppsx
<input checked="" type="checkbox"/>	docx
<input checked="" type="checkbox"/>	pptx

Torrent Metadata Parser

<input checked="" type="checkbox"/>	torrent
-------------------------------------	---------

Word Document Parser

<input checked="" type="checkbox"/>	doc
-------------------------------------	-----

OASIS OpenDocument V2 Text Document

<input checked="" type="checkbox"/>	odg
<input checked="" type="checkbox"/>	odf

Bzip 2 UNIX Compressed File Parser

<input checked="" type="checkbox"/>	tbz
<input checked="" type="checkbox"/>	tbz2
<input checked="" type="checkbox"/>	bz2

Streaming HTML Parser

<input checked="" type="checkbox"/>	xhtml
<input checked="" type="checkbox"/>	php4
<input checked="" type="checkbox"/>	php5
<input checked="" type="checkbox"/>	php3
<input checked="" type="checkbox"/>	shtml
<input checked="" type="checkbox"/>	html
<input checked="" type="checkbox"/>	htm

Microsoft Excel Parser

<input checked="" type="checkbox"/>	xla
<input checked="" type="checkbox"/>	xls

ZIP File Parser

<input checked="" type="checkbox"/>	zip
<input checked="" type="checkbox"/>	jar
<input checked="" type="checkbox"/>	apk

Acrobat Portable Document Parser

<input checked="" type="checkbox"/>	pdf
-------------------------------------	-----

7zip Archive Parser

<input checked="" type="checkbox"/>	7z
-------------------------------------	----

RSS Parser

<input checked="" type="checkbox"/>	rss
<input checked="" type="checkbox"/>	xml

Tape Archive File Parser

<input checked="" type="checkbox"/>	tar
-------------------------------------	-----

Rich Text Format Parser

<input checked="" type="checkbox"/>	rtf
-------------------------------------	-----

RDF Parser

<input checked="" type="checkbox"/>	rdf
-------------------------------------	-----

Submit

Komponenten der Suchmaschine

crawler

parser

Index Schema



Active	Attribute	Comment	show active	show all available	show disabled
<input checked="" type="checkbox"/>	id	primary key of document, the URL hash mandatory field			
<input checked="" type="checkbox"/>	sku	url of document			
<input checked="" type="checkbox"/>	last_modified	last-modified from http header			
<input type="checkbox"/>	dates_in_content_sxt	if date expressions can be found in the content, these dates are listed here in order of the appearances			
<input type="checkbox"/>	dates_in_content_count_i	the number of entries in dates_in_content_sxt			
<input type="checkbox"/>	date_in_content_min_dt	if dates_in_content_sxt is filled, this contains the oldest date from the list of available dates			
<input type="checkbox"/>	date_in_content_max_dt	if dates_in_content_sxt is filled, this contains the youngest date from the list of available dates, that may also be possibly in the future			
<input checked="" type="checkbox"/>	content_type	mime-type of document			
<input type="checkbox"/>	http_unique_b	unique-field which is true when an url appears the first time. If the same url which was http then appears as https (or vice versa) then the field is false			
<input type="checkbox"/>	www_unique_b	unique-field which is true when an url appears the first time. If the same url within the subdomain www then appears without that subdomain (or vice versa) then the field is false			
<input checked="" type="checkbox"/>	title	content of title tag			
<input type="checkbox"/>	title_exact_signature_l	the 64 bit hash of the org.apache.solr.update.processor.Lookup3Signature of title, used to compute title_unique_b			
<input type="checkbox"/>	title_unique_b	flag shows if title is unique within all indexable documents of the same host with status code 200; if yes and another document appears with same title, the unique-flag is set to false			
<input checked="" type="checkbox"/>	host_id_s	id of the host, a 6-byte hash that is part of the document id			
<input checked="" type="checkbox"/>	md5_s	the md5 of the raw source			
<input checked="" type="checkbox"/>	exact_signature_l	the 64 bit hash of the org.apache.solr.update.processor.Lookup3Signature of text_t			
<input checked="" type="checkbox"/>	exact_signature_unique_b	flag shows if exact_signature_l is unique at the time of document creation, used for double-check during search			
<input type="checkbox"/>	exact_signature_copycount_i	counter for the number of documents which are not unique (== count of not-unique-flagged documents + 1)			
<input checked="" type="checkbox"/>	fuzzy_signature_l	64 bit of the Lookup3Signature from EnhancedTextProfileSignature of text_t			
<input type="checkbox"/>	fuzzy_signature_text_t	intermediate data produced in EnhancedTextProfileSignature: a list of word frequencies			
<input checked="" type="checkbox"/>	fuzzy_signature_unique_b	flag shows if fuzzy_signature_l is unique at the time of document creation, used for double-check during search			
<input type="checkbox"/>	fuzzy_signature_copycount_i	counter for the number of documents which are not unique (== count of not-unique-flagged documents + 1)			
<input checked="" type="checkbox"/>	size_i	the size of the raw source			
<input checked="" type="checkbox"/>	failreason_s	fail reason if a page was not loaded. if the page was loaded then this field is empty			
<input checked="" type="checkbox"/>	failtype_s	fail type if a page was not loaded. This field is either empty, 'excl' or 'fail'			
<input checked="" type="checkbox"/>	httpstatus_i	html status return code (i.e. "200" for ok), -1 if not loaded			
<input type="checkbox"/>	httpstatus_redirect_s	redirect url if the error code is 299 < httpstatus_i < 310			
<input checked="" type="checkbox"/>	references_i	number of unique http references, should be equal to references_internal_i + references_external_i			
<input checked="" type="checkbox"/>	references_internal_i	number of unique http references from same host to referenced url			
<input checked="" type="checkbox"/>	references_external_i	number of unique http references from external hosts			
<input checked="" type="checkbox"/>	references_exthosts_i	number of external hosts which provide http references			
<input type="checkbox"/>	crawldepth_i	crawl depth of web page according to the number of steps that the crawler did to get to this document; if the crawl was started at a root			

Funktionskomponente

Komponenten der Suchmaschine

crawler

parser

Index Schema



Information Retrieval mit YaCy

<input type="checkbox"/>	httpstatus_redirect_s	redirect url if the error code is 299 < httpstatus_j < 310
<input checked="" type="checkbox"/>	references_j	number of unique http references, should be equal to references_internal_j + references_external_j
<input checked="" type="checkbox"/>	references_internal_j	number of unique http references from same host to referenced url
<input checked="" type="checkbox"/>	references_external_j	number of unique http references from external hosts
<input checked="" type="checkbox"/>	references_exthosts_j	number of external hosts which provide http references
<input checked="" type="checkbox"/>	crawlddepth_j	crawl depth of web page according to the number of steps that the crawler did to get to this document; if the crawl was started at a root document, then this is equal to the clickdepth
<input checked="" type="checkbox"/>	process_sxt	needed (post-)processing steps on this metadata set
<input checked="" type="checkbox"/>	harvestkey_s	key from a harvest process (i.e. the crawl profile hash key) which is needed for near-realtime postprocessing. This shall be deleted as soon as postprocessing has been terminated.
<input checked="" type="checkbox"/>	load_date_dt	time when resource was loaded
<input checked="" type="checkbox"/>	fresh_date_dt	date until resource shall be considered as fresh
<input checked="" type="checkbox"/>	referrer_id_s	id of the referrer to this document, discovered during crawling
<input checked="" type="checkbox"/>	publisher_t	the name of the publisher of the document
<input checked="" type="checkbox"/>	language_s	the language used in the document
<input checked="" type="checkbox"/>	audiolinkscount_j	number of links to audio resources
<input checked="" type="checkbox"/>	videolinkscount_j	number of links to video resources
<input checked="" type="checkbox"/>	applinkscount_j	number of links to application resources
<input checked="" type="checkbox"/>	coordinate_p	point in degrees of latitude,longitude as declared in WSG84
<input type="checkbox"/>	coordinate_p_0_coordinate	automatically created subfield, (latitude)
<input type="checkbox"/>	coordinate_p_1_coordinate	automatically created subfield, (longitude)
<input type="checkbox"/>	ip_s	ip of host of url (after DNS lookup)
<input checked="" type="checkbox"/>	author	content of author-tag
<input type="checkbox"/>	author_sxt	content of author-tag as copy-field from author. This is used for facet generation
<input checked="" type="checkbox"/>	description_bxt	content of description-tag(s)
<input type="checkbox"/>	description_exact_signature_j	the 64 bit hash of the org.apache.solr.update.processor.Lookup3Signature of description, used to compute description_unique_b
<input type="checkbox"/>	description_unique_b	flag shows if description is unique within all indexable documents of the same host with status code 200; if yes and another document appears with same description, the unique-flag is set to false
<input checked="" type="checkbox"/>	keywords	content of keywords tag; words are separated by space
<input checked="" type="checkbox"/>	charset_s	character encoding
<input checked="" type="checkbox"/>	wordcount_j	number of words in visible area
<input checked="" type="checkbox"/>	linkscout_j	number of all outgoing links; including linksnofollowcount_j
<input checked="" type="checkbox"/>	linksnofollowcount_j	number of all outgoing inks with nofollow tag
<input checked="" type="checkbox"/>	inboundlinkscout_j	number of outgoing inbound (to same domain) links; including inboundlinksnofollowcount_j
<input checked="" type="checkbox"/>	inboundlinksnofollowcount_j	number of outgoing inbound (to same domain) links with nofollow tag
<input checked="" type="checkbox"/>	outboundlinkscout_j	number of outgoing outbound (to other domain) links, including outboundlinksnofollowcount_j
<input checked="" type="checkbox"/>	outboundlinksnofollowcount_j	number of outgoing outbound (to other domain) links with nofollow tag
<input checked="" type="checkbox"/>	imagescount_j	number of images
<input checked="" type="checkbox"/>	responsetime_j	response time of target server in milliseconds
<input checked="" type="checkbox"/>	text_t	all visible text
<input checked="" type="checkbox"/>	synonyms_sxt	additional synonyms to the words in the text
<input type="checkbox"/>	h1_bxt	h1 header

Funktionskomponenten

Komponenten der Suchmaschine

crawler

parser

Index Schema



Information Retrieval mit YaCy

<input checked="" type="checkbox"/>	inboundlinks_nofollowcount_i	number of outgoing inbound (to same domain) links with nofollow tag
<input checked="" type="checkbox"/>	outboundlinkscount_i	number of outgoing outbound (to other domain) links, including outboundlinksnofollowcount_i
<input checked="" type="checkbox"/>	outboundlinksnofollowcount_i	number of outgoing outbound (to other domain) links with nofollow tag
<input checked="" type="checkbox"/>	imagescount_i	number of images
<input checked="" type="checkbox"/>	responsetime_i	response time of target server in milliseconds
<input checked="" type="checkbox"/>	text_t	all visible text
<input checked="" type="checkbox"/>	synonyms_sxt	additional synonyms to the words in the text
<input checked="" type="checkbox"/>	h1_txt	h1 header
<input checked="" type="checkbox"/>	h2_txt	h2 header
<input checked="" type="checkbox"/>	h3_txt	h3 header
<input checked="" type="checkbox"/>	h4_txt	h4 header
<input checked="" type="checkbox"/>	h5_txt	h5 header
<input checked="" type="checkbox"/>	h6_txt	h6 header
<input checked="" type="checkbox"/>	collection_sxt	tags that are attached to crawls/index generation to separate the search result into user-defined subsets
<input type="checkbox"/>	csscount_i	number of entries in css_tag_sxt and css_url_sxt
<input type="checkbox"/>	css_tag_sxt	full css tag with normalized url
<input type="checkbox"/>	css_url_sxt	normalized urls within a css tag
<input type="checkbox"/>	scripts_sxt	normalized urls within a scripts tag
<input type="checkbox"/>	scriptscount_i	number of entries in scripts_sxt
<input type="checkbox"/>	robots_i	content of <meta name="robots" content=#content#> tag and the "X-Robots-Tag" HTTP property
<input type="checkbox"/>	metagenerator_t	content of <meta name="generator" content=#content#> tag
<input checked="" type="checkbox"/>	inboundlinks_protocol_sxt	internal links, only the protocol
<input checked="" type="checkbox"/>	inboundlinks_urlstub_sxt	internal links, the url only without the protocol
<input checked="" type="checkbox"/>	inboundlinks_anchortext_txt	internal links, the visible anchor text
<input checked="" type="checkbox"/>	outboundlinks_protocol_sxt	external links, only the protocol
<input checked="" type="checkbox"/>	outboundlinks_urlstub_sxt	external links, the url only without the protocol
<input checked="" type="checkbox"/>	outboundlinks_anchortext_txt	external links, the visible anchor text
<input checked="" type="checkbox"/>	images_text_t	all text/words appearing in image alt texts or the tokenized url
<input checked="" type="checkbox"/>	images_urlstub_sxt	all image links without the protocol and '://'
<input checked="" type="checkbox"/>	images_protocol_sxt	all image link protocols
<input checked="" type="checkbox"/>	images_alt_sxt	all image link alt tag
<input checked="" type="checkbox"/>	images_height_val	size of images:height
<input checked="" type="checkbox"/>	images_width_val	size of images:width
<input type="checkbox"/>	images_pixel_val	size of images as number of pixels (easier for a search restriction than with and height)
<input type="checkbox"/>	images_withalt_i	number of image links with alt tag
<input type="checkbox"/>	htags_i	binary pattern for the existence of h1..h6 headlines
<input type="checkbox"/>	canonical_s	url inside the canonical link element
<input type="checkbox"/>	canonical_equal_sku_b	flag shows if the url in canonical_t is equal to sku
<input type="checkbox"/>	refresh_s	link from the url property inside the refresh link element
<input type="checkbox"/>	li_txt	all texts in tags

Funktionskomponenten

Komponenten der Suchmaschine

crawler

parser

Index Schema



Information Retrieval mit YaCy

<input type="checkbox"/>	canonical_s	url inside the canonical link element
<input type="checkbox"/>	canonical_equal_sku_b	flag shows if the url in canonical_t is equal to sku
<input type="checkbox"/>	refresh_s	link from the url property inside the refresh link element
<input type="checkbox"/>	li_txt	all texts in tags
<input type="checkbox"/>	licount_i	number of tags
<input checked="" type="checkbox"/>	bold_txt	all texts inside of or tags. no doubles. listed in the order of number of occurrences in decreasing order
<input type="checkbox"/>	boldcount_i	total number of occurrences of or
<input checked="" type="checkbox"/>	italic_txt	all texts inside of <i> tags. no doubles. listed in the order of number of occurrences in decreasing order
<input type="checkbox"/>	italiccount_i	total number of occurrences of <i>
<input checked="" type="checkbox"/>	underline_txt	all texts inside of <u> tags. no doubles. listed in the order of number of occurrences in decreasing order
<input type="checkbox"/>	underlinecount_i	total number of occurrences of <u>
<input type="checkbox"/>	flash_b	flag that shows if a swf file is linked
<input type="checkbox"/>	frames_sxt	list of all links to frames
<input type="checkbox"/>	framesscount_i	number of frames_txt
<input type="checkbox"/>	iframes_sxt	list of all links to iframes
<input type="checkbox"/>	iframesscount_i	number of iframes_txt
<input type="checkbox"/>	hreflang_url_sxt	url of the hreflang link tag, see http://support.google.com/webmasters/bin/answer.py?hl=de&answer=189077
<input type="checkbox"/>	hreflang_cc_sxt	country code of the hreflang link tag, see http://support.google.com/webmasters/bin/answer.py?hl=de&answer=189077
<input type="checkbox"/>	navigation_url_sxt	page navigation url, see http://googlewebmastercentral.blogspot.de/2011/09/pagination-with-relnext-and-relprev.html
<input type="checkbox"/>	navigation_type_sxt	page navigation rel property value, can contain one of {top,up,next,prev,first,last}
<input type="checkbox"/>	publisher_url_s	publisher url as defined in http://support.google.com/plus/answer/1713826?hl=de
<input checked="" type="checkbox"/>	url_protocol_s	the protocol of the url
<input checked="" type="checkbox"/>	url_file_name_s	the file name (which is the string after the last '/' and before the query part from '?' on) without the file extension
<input type="checkbox"/>	url_file_name_tokens_t	tokens generated from url_file_name_s which can be used for better matching and result boosting
<input checked="" type="checkbox"/>	url_file_ext_s	the file name extension
<input checked="" type="checkbox"/>	url_paths_count_i	number of all path elements in the url hpath (see: http://www.ietf.org/rfc/rfc1738.txt) without the file name
<input checked="" type="checkbox"/>	url_paths_sxt	all path elements in the url hpath (see: http://www.ietf.org/rfc/rfc1738.txt) without the file name
<input type="checkbox"/>	url_parameter_i	number of key-value pairs in search part of the url
<input type="checkbox"/>	url_parameter_key_sxt	the keys from key-value pairs in the search part of the url
<input type="checkbox"/>	url_parameter_value_sxt	the values from key-value pairs in the search part of the url
<input checked="" type="checkbox"/>	url_chars_i	number of all characters in the url == length of sku field
<input checked="" type="checkbox"/>	host_s	host of the url
<input type="checkbox"/>	host_dnc_s	the Domain Class Name, either the TLD or a combination of ccSLD+TLD if a ccSLD is used.
<input checked="" type="checkbox"/>	host_organization_s	either the second level domain or, if a ccSLD is used, the third level domain
<input type="checkbox"/>	host_organizationdnc_s	the organization and dnc concatenated with '.'
<input type="checkbox"/>	host_subdomain_s	the remaining part of the host without organizationdnc
<input checked="" type="checkbox"/>	host_extent_i	number of documents from the same host; can be used to measure references_internal_i for likelihood computation
<input type="checkbox"/>	title_count_i	number of titles (counting the 'title' field) in the document
<input type="checkbox"/>	title_chars_val	number of characters for each title

Funktionskomponenten

Komponenten der Suchmaschine

crawler

parser

Index Schema



<input type="checkbox"/>	title_chars_val	number of characters for each title
<input type="checkbox"/>	title_words_val	number of words in each title
<input type="checkbox"/>	description_count_i	number of descriptions in the document. Its not counting the 'description' field since there is only one. But it counts the number of descriptions that appear in the document (if any)
<input type="checkbox"/>	description_chars_val	number of characters for each description
<input type="checkbox"/>	description_words_val	number of words in each description
<input type="checkbox"/>	h1_i	number of h1 header lines
<input type="checkbox"/>	h2_i	number of h2 header lines
<input type="checkbox"/>	h3_i	number of h3 header lines
<input type="checkbox"/>	h4_i	number of h4 header lines
<input type="checkbox"/>	h5_i	number of h5 header lines
<input type="checkbox"/>	h6_i	number of h6 header lines
<input type="checkbox"/>	schema_org_breadcrumb_i	number of itemprop="breadcrumb" appearances in div tags
<input type="checkbox"/>	opengraph_title_t	Open Graph Metadata from og:title metadata field, see http://ogp.me/ns#
<input type="checkbox"/>	opengraph_type_s	Open Graph Metadata from og:type metadata field, see http://ogp.me/ns#
<input type="checkbox"/>	opengraph_url_s	Open Graph Metadata from og:url metadata field, see http://ogp.me/ns#
<input type="checkbox"/>	opengraph_image_s	Open Graph Metadata from og:image metadata field, see http://ogp.me/ns#
<input type="checkbox"/>	cr_host_count_i	the number of documents within a single host
<input type="checkbox"/>	cr_host_chance_d	the chance to click on this page when randomly clicking on links within on one host
<input type="checkbox"/>	cr_host_norm_i	normalization of chance: 0 for lower half of cr_host_count_i urls, 1 for 1/2 of the remaining and so on. the maximum number is 10
<input type="checkbox"/>	rating_i	custom rating; to be set with external rating information
<input type="checkbox"/>	bold_val	number of occurrences of texts in bold_txt
<input type="checkbox"/>	italic_val	number of occurrences of texts in italic_txt
<input type="checkbox"/>	underline_val	number of occurrences of texts in underline_txt
<input type="checkbox"/>	ext_cms_txt	names of cms attributes; if several are recognized then they are listen in decreasing order of number of matching criterias
<input type="checkbox"/>	ext_cms_val	number of attributes that count for a specific cms in ext_cms_txt
<input type="checkbox"/>	ext_ads_txt	names of ad-servers/ad-services
<input type="checkbox"/>	ext_ads_val	number of attributes counts in ext_ads_txt
<input type="checkbox"/>	ext_community_txt	names of recognized community functions
<input type="checkbox"/>	ext_community_val	number of attribute counts in attr_community
<input type="checkbox"/>	ext_maps_txt	names of map services
<input type="checkbox"/>	ext_maps_val	number of attribute counts in ext_maps_txt
<input type="checkbox"/>	ext_tracker_txt	names of tracker server
<input type="checkbox"/>	ext_tracker_val	number of attribute counts in ext_tracker_txt
<input type="checkbox"/>	ext_title_txt	names matching title expressions
<input type="checkbox"/>	ext_title_val	number of matching title expressions
<input checked="" type="checkbox"/>	vocabularies_sxt	collection of all vocabulary names that have a matcher in the document - use this to boost with vocabularies

Funktionskomponenten einer Suchmaschine



Textfelder aus Index Schema

sku	<input type="checkbox"/>		url of document
title	<input checked="" type="checkbox"/>	5.0	content of title tag
fuzzy_signature_text_t	<input type="checkbox"/>		intermediate data produced in
publisher_t	<input type="checkbox"/>		the name of the publisher of the document
author	<input type="checkbox"/>		content of author-tag
description_txt	<input type="checkbox"/>		content of description-tag(s)
keywords	<input type="checkbox"/>		content of keywords tag; words are separated by
text_t	<input checked="" type="checkbox"/>	1.0	all visible text
synonyms_sxt	<input checked="" type="checkbox"/>	0.5	additional synonyms to the words in the text
h1_txt	<input checked="" type="checkbox"/>	5.0	h1 header
h2_txt	<input checked="" type="checkbox"/>	3.0	h2 header
h3_txt	<input type="checkbox"/>		h3 header
h4_txt	<input type="checkbox"/>		h4 header
h5_txt	<input type="checkbox"/>		h5 header
h6_txt	<input type="checkbox"/>		h6 header
inboundlinks_urlstub_sxt	<input type="checkbox"/>		internal links, th
inboundlinks_anchor_text_txt	<input type="checkbox"/>		internal links, th
outboundlinks_urlstub_sxt	<input type="checkbox"/>		external links, th

outboundlinks_anchor_text_txt	<input type="checkbox"/>		external links, the visible anchor text
images_text_t	<input type="checkbox"/>		all text/words appearing in image alt texts or the
images_urlstub_sxt	<input type="checkbox"/>		all image links without the protocol and '://'
images_alt_sxt	<input type="checkbox"/>		all image link alt tag
li_txt	<input type="checkbox"/>		all texts in tags
bold_txt	<input type="checkbox"/>		all texts inside of or tags. no doubles.
italic_txt	<input type="checkbox"/>		all texts inside of <i> tags. no doubles. listed in the
underline_txt	<input type="checkbox"/>		all texts inside of <u> tags. no doubles. listed in the
url_file_name_s	<input type="checkbox"/>		the file name (which is the string after the last '/')
url_file_name_tokens_t	<input checked="" type="checkbox"/>	4.0	tokens generated from url_file_name_s which can
url_file_ext_s	<input type="checkbox"/>		the file name extension
url_paths_sxt	<input checked="" type="checkbox"/>	3.0	all path elements in the url hpath (see:
host_s	<input checked="" type="checkbox"/>	6.0	host of the url
host_dnc_s	<input type="checkbox"/>		the Domain Class Name, either the TLD or a
host_organization_s	<input type="checkbox"/>		either the second level domain or, if a ccSLD is
host_organizationdnc_s	<input type="checkbox"/>		the organization and dnc concatenated with '.'
host_subdomain_s	<input type="checkbox"/>		the remaining part of the host without
opengraph_title_t	<input type="checkbox"/>		Open Graph Metadata from og:title metadata field,



Funktionskomponenten einer Suchmaschine



http://localhost:8090/solr/select?q=text_t:ibm%20mainframe%20AND%20url_file_ext_s:pdf&fl=sku,author,publisher_t

Entwicklung einer eigenen Suchmaschine

crawler

parser

search interface

Apache
Solr



Web Search



ibm mainframe filetype:pdf

3 results from a total of 3 docs in index; search time: 217 milliseconds.

[create a download script](#) [all results](#)

[remove the filter 'filetype:pdf'](#)

Count	Protocol	Host	Path	URL	Size	Date
1	https	www.dbsystel.de	/file/3275180/data/	https://www.dbsystel.de/file/3275180/data/bereitstellung_betrieb_von_ibm-mainframe-kapazitaeten_englisch.pdf	2 mbyte	Mon, 05 Aug 2013 08:14:04
2	https	www.dbsystel.de	/file/2239284/data/	https://www.dbsystel.de/file/2239284/data/bereitstellung_betrieb_von_ibm-mainframe-kapazitaeten_deutsch.pdf	2 mbyte	Wed, 16 Mar 2011 16:36:45
3	https	www.dbsystel.de	/file/2247366/data/	https://www.dbsystel.de/file/2247366/data/bereitstellung_betrieb_von_ibm-mainframe-kapazitaeten_englisch.pdf	2 mbyte	Tue, 04 Oct 2011 07:04:25

Based Trai [...]

<http://ma-co.de/>

Sat, 21 Sep 2013 | [Metadata](#) | [Parser](#) | [Citations](#) | [*](#)

Management- und Personalberatung van der Zalm Hamburg

[...] Suche Direktansprache MANAGEMENT Potenzialanalyse Organisationsentwicklung **Personalentwicklung** Vergütungsmanagement TRAINING & COACHING. [...]

<http://www.vanderzalm.de/index.php?kat=997>

Sat, 21 Sep 2013 | [Metadata](#) | [Parser](#) | [Citations](#) | [*](#)

Preise und Auszeichnungen für EF Englishtown

Jedes Jahr werden ausgewählte Personalabteilungen und Trainingsanbieter vom HRE für ihr herausragendes Personalmanagement und vorbildliche **Personalentwicklung** geehrt

<http://www.englishtown.de/online/awards.aspx>

Sat, 21 Sep 2013 | [Metadata](#) | [Parser](#) | [Citations](#) | [*](#)

mpmEXPERTS - Wir machen MultiProjektmanagement. Einfach. Sicher.

Personalentwicklung, Talanx Service AG & 100% professionell, hervorragende

</result>

</response>

- dbsystel.de (9)
- stadtbranche.de (5)
- mpm-experts.com (2)
- d-nb.info (1)
- berlin.city-map.de (1)
- englishtown.de (1)
- de.wikipedia.org (1)
- lebensmittel-verzeichnis.de (1)
- central.de (1)
- nammert24.de (1)
- ma-co.de (1)
- vanderzalm.de (1)

Studies Navigator

- Seminars (1)



Komponenten der Suchmaschine

crawler

parser

search interface

monitoring

administration



Process Scheduler

Recorded Actions

<input type="checkbox"/>	Type	Comment	Call Count	Recording Date	Last Exec Date	Next Exec Date	Event Trigger	Scheduler
<input type="checkbox"/>	crawler	crawl start for http://www.geschichteinchronologie.ch/	1	Sep 2, 2013 7:49:21 PM	Sep 2, 2013 7:49:21 PM	-	no event	no repetition
<input type="checkbox"/>	crawler	crawl start for http://foldoc.org/contents/all.html	1	Oct 25, 2013 7:31:46 PM	Oct 25, 2013 7:31:46 PM	-	no event	no repetition
<input type="checkbox"/>	crawler	crawl start for http://www.heinrich-kromer-schule.de/	1	Nov 3, 2013 12:55:37 PM	Nov 3, 2013 12:55:37 PM	-	no event	no repetition
<input type="checkbox"/>	crawler	crawl start for http://www.dbsystel.de/dbsystel/start/	1	Nov 3, 2013 7:18:07 PM	Nov 3, 2013 7:18:07 PM	-	no event	7 days
<input type="checkbox"/>	crawler	crawl start for http://www.cafe-nuesslein.de/	1	Nov 3, 2013 7:19:07 PM	Nov 3, 2013 7:19:07 PM	-	run regular	no repetition

Execute Selected Actions

Delete Selected Actions

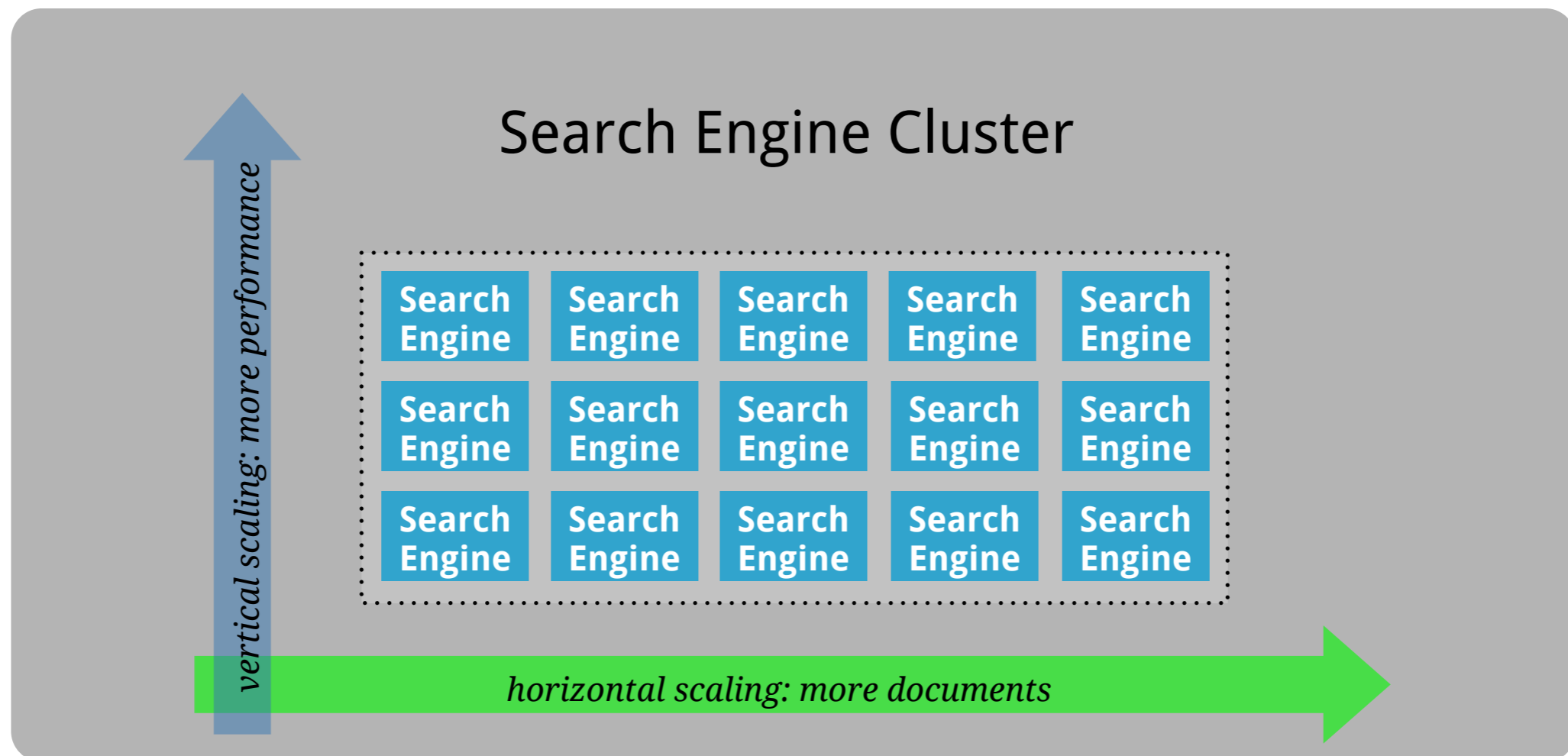
Delete all Actions which had been created before

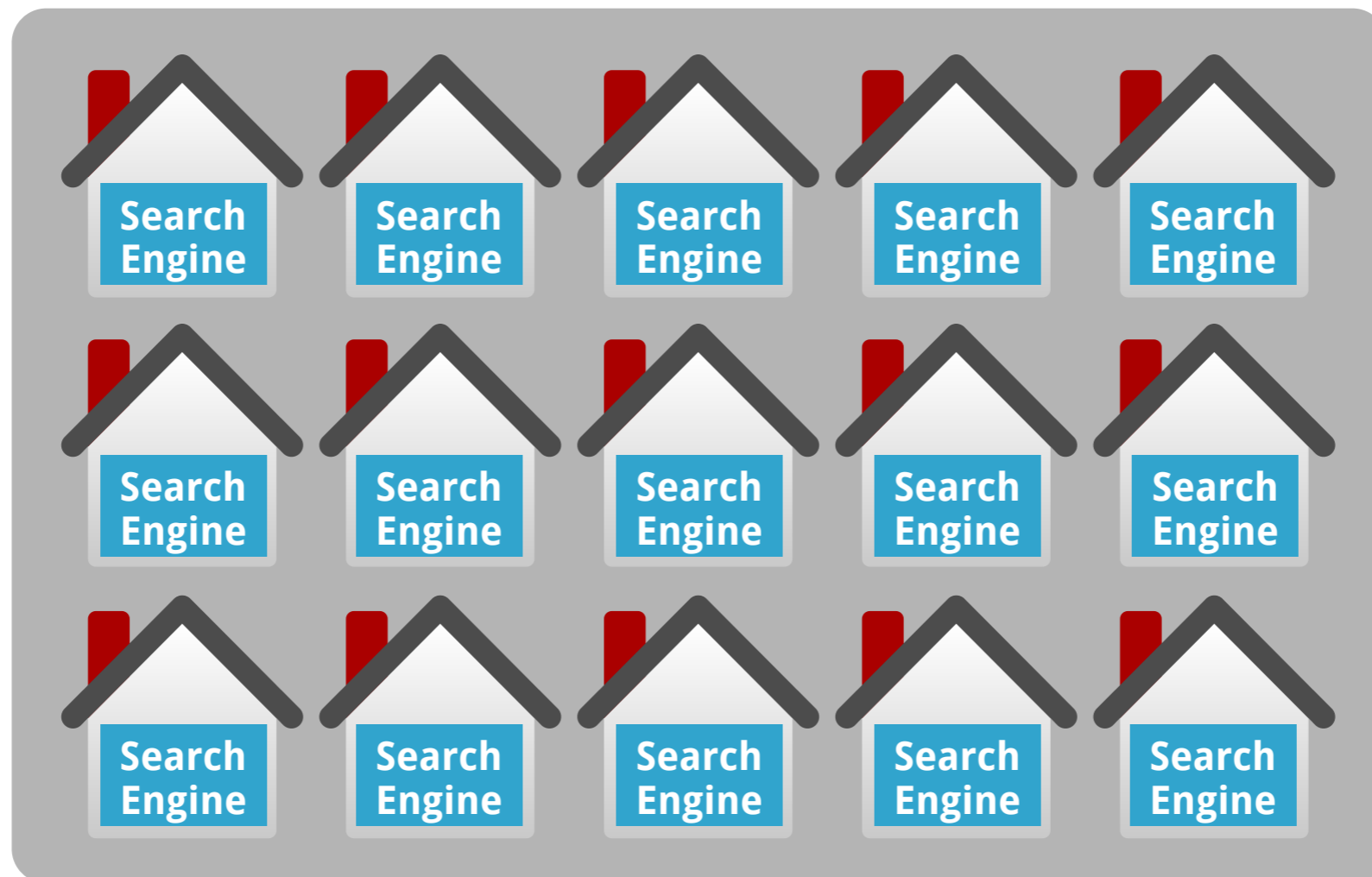
2 months

no event 7 days

run regular no repetition

after start-up





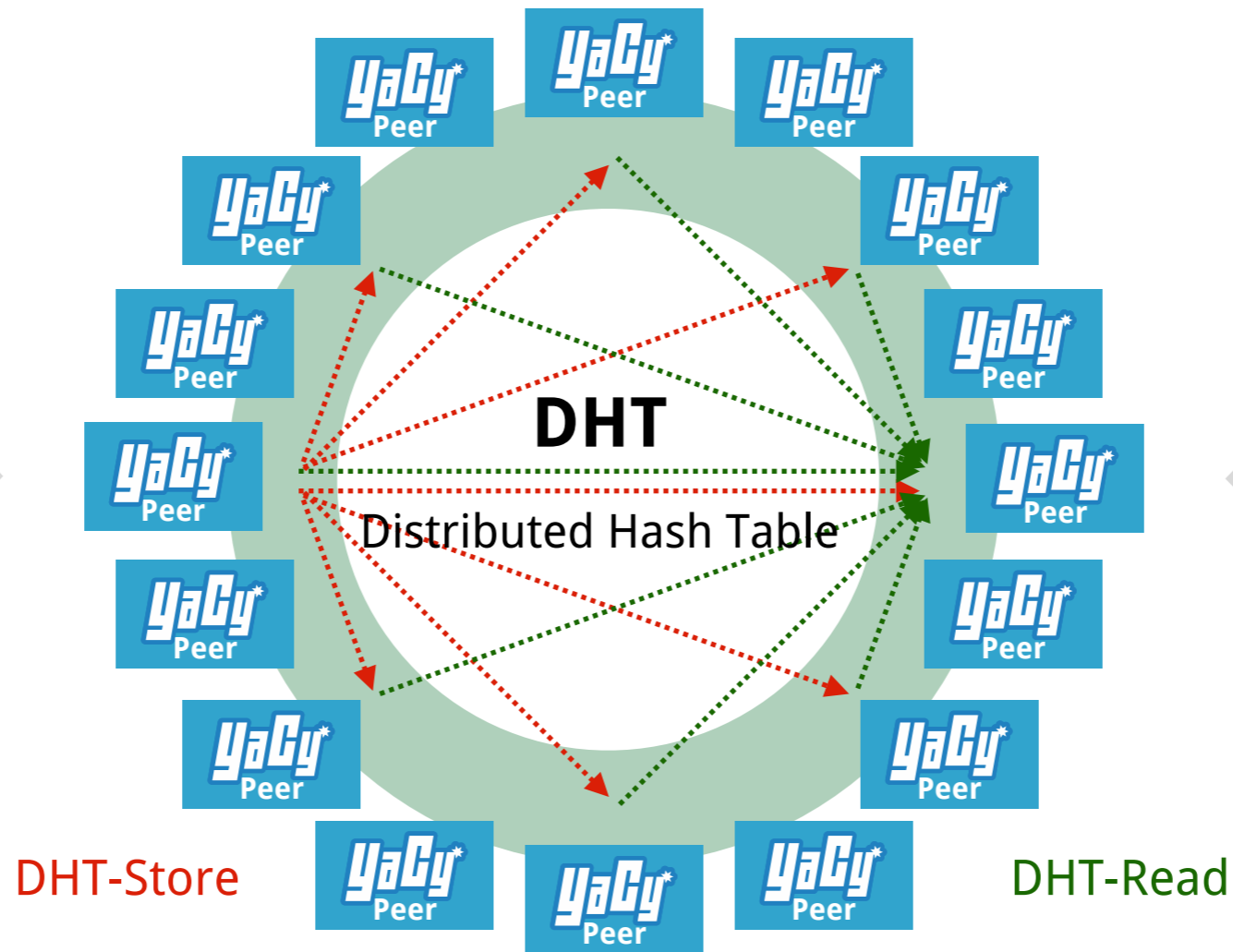
Skalierung Distributed + Decentralized



Skalierung Peer-to-Peer Distributed + Decentralized



Crawl the web, create a web index, distribute the index



Search in a Distributed Hash Table

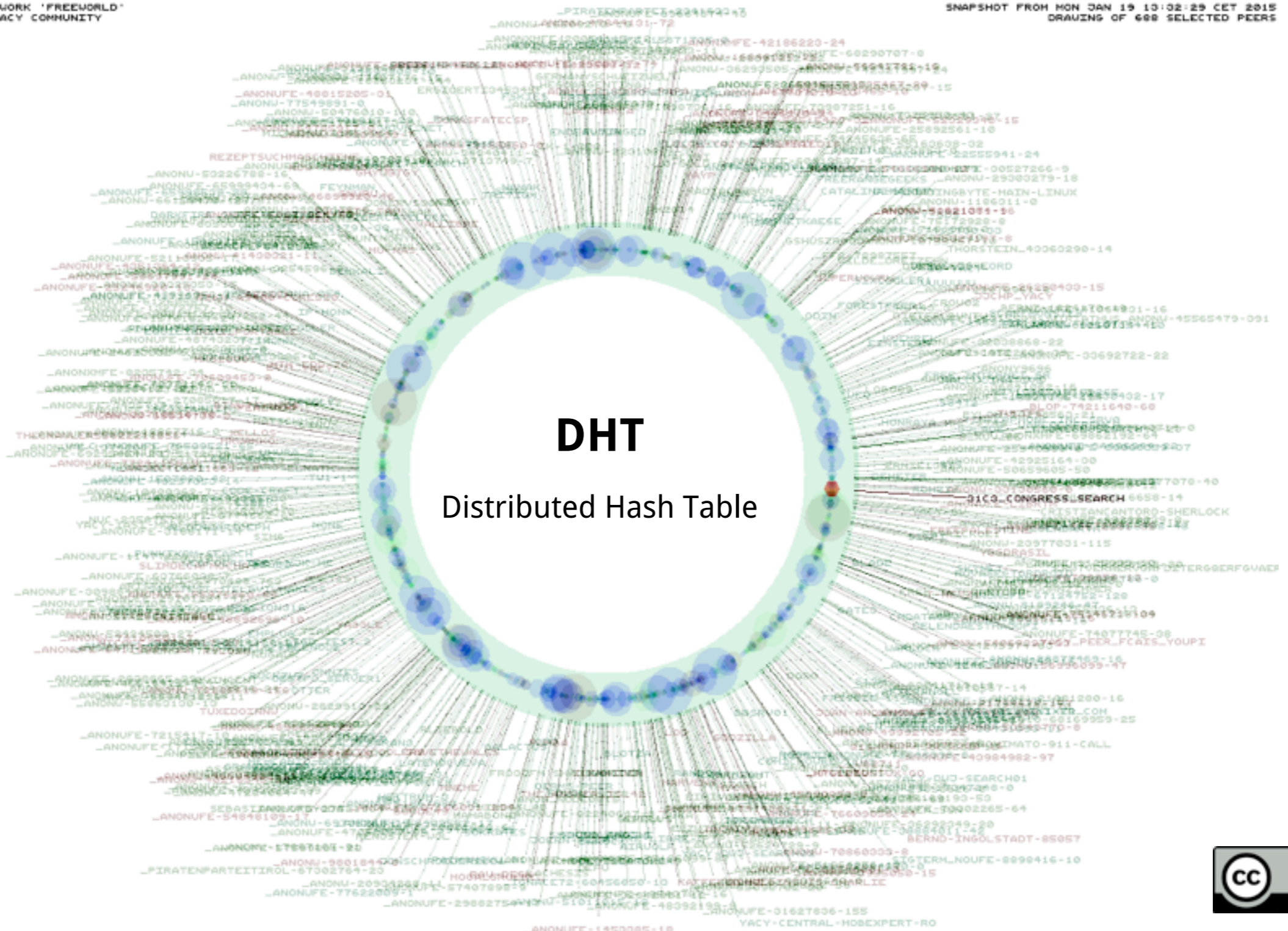


Skalierung Peer-to-Peer Distributed + Decentralized



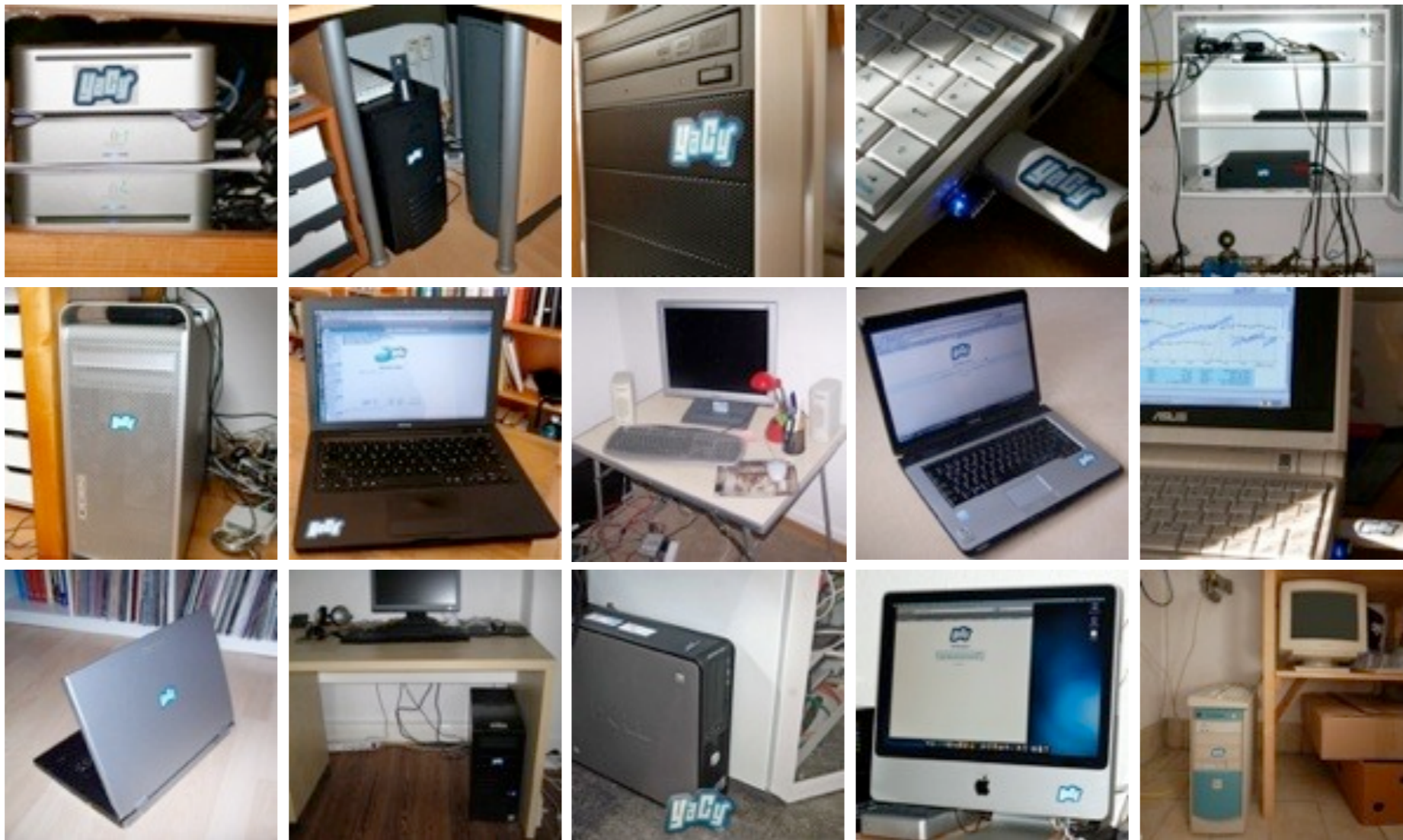
YACY NETWORK 'FREEWORLD'
PUBLIC YACY COMMUNITY

SNAPSHOT FROM MON JAN 19 10:02:29 CET 2015
DRAWING OF 600 SELECTED PEERS



Das YaCy ‚freeworld‘ Netz

> 1 Milliarde Dokumente, ca. 500 Nutzer/Tag




Solr Boosts

This is the set of searchable fields (see [YaCy Solr Schema](#)). Entries without a boost value are not searched. Boost values make hits inside the corresponding field more important.

sku	<input type="checkbox"/>		url of document
title	<input checked="" type="checkbox"/>	5.0	content of title tag
publisher_t	<input type="checkbox"/>		the name of the publisher of the document
author	<input type="checkbox"/>		content of author-tag
description_txt	<input type="checkbox"/>		content of description-tag(s)
keywords	<input type="checkbox"/>		content of keywords tag; words are separated by
text_t	<input checked="" type="checkbox"/>	1.0	all visible text
synonyms_sxt	<input checked="" type="checkbox"/>	0.5	additional synonyms to the words in the text
h1_txt	<input checked="" type="checkbox"/>	5.0	h1 header
h2_txt	<input checked="" type="checkbox"/>	3.0	h2 header
h3_txt	<input type="checkbox"/>		h3 header
h4_txt	<input type="checkbox"/>		h4 header
h5_txt	<input type="checkbox"/>		h5 header
h6_txt	<input type="checkbox"/>		h6 header
inboundlinks_urlstub_sxt	<input type="checkbox"/>		internal links, the url only without the protocol
inboundlinks_anchor_text_txt	<input type="checkbox"/>		internal links, the visible anchor text
outboundlinks_urlstub_sxt	<input type="checkbox"/>		external links, the url only without the protocol
outboundlinks_anchor_text_txt	<input type="checkbox"/>		external links, the visible anchor text
images_text_t	<input type="checkbox"/>		all text/words appearing in image alt texts or the
images_urlstub_sxt	<input type="checkbox"/>		all image links without the protocol and '//'
images_alt_sxt	<input type="checkbox"/>		all image link alt tag
bold_txt	<input type="checkbox"/>		all texts inside of or tags. no doubles.
italic_txt	<input type="checkbox"/>		all texts inside of <i> tags. no doubles. listed in the
underline_txt	<input type="checkbox"/>		all texts inside of <u> tags. no doubles. listed in the
url_file_name_s	<input type="checkbox"/>		the file name (which is the string after the last '/')
url_file_name_tokens_t	<input checked="" type="checkbox"/>	4.0	field not in local index (boost has no effect) tokens
url_file_ext_s	<input type="checkbox"/>		the file name extension
url_paths_sxt	<input checked="" type="checkbox"/>	3.0	all path elements in the url hpath (see:
host_s	<input checked="" type="checkbox"/>	6.0	host of the url
host_organization_s	<input type="checkbox"/>		either the second level domain or, if a ccSLD is

[Set Field Boosts](#) [Re-Set to default](#)

Boost Function

A Boost Function can combine numeric values from the result document to produce a number which is multiplied with the score value from the query result. To see all available fields, see the [YaCy Solr Schema](#) and look for numeric values (these are names with suffix '_i'). To find out which kind of operations are possible, see the [Solr Function Query](#)  documentation. Example: to order by date, use "recip(ms(NOW,last_modified),3.16e-11,1,1)", to order by crawldepth, use "div(100,add(crawldepth_i,1))".

boost=

Set Boost Function

Re-Set to default

You can boost with vocabularies, use the occurrence counters [vocabulary_Dokumentenart_i, vocabulary_Locations_i] and [vocabulary_Dokumentenart_log_i, vocabulary_Locations_log_i].

Boost Query

The Boost Query is attached to every query. Use this to statically boost specific content in the index. Example: "fuzzy_signature_unique_b:true^100000.0f" means that documents, identified as 'double' are ranked very bad and appended to the end of all results (because the unique are ranked high). To find appropriate fields for this query, see the [YaCy Solr Schema](#) and look for boolean values (with suffix '_b') or tags inside string fields (with suffix '_s' or '_sxt').

bq=crawldepth_i:0^0.8 crawldepth_i:1^0.4 url_protocol_s:https^10.0 http_unique_b:[* TO *]^100.0

Set Boost Query

Re-Set to default

Filter Query

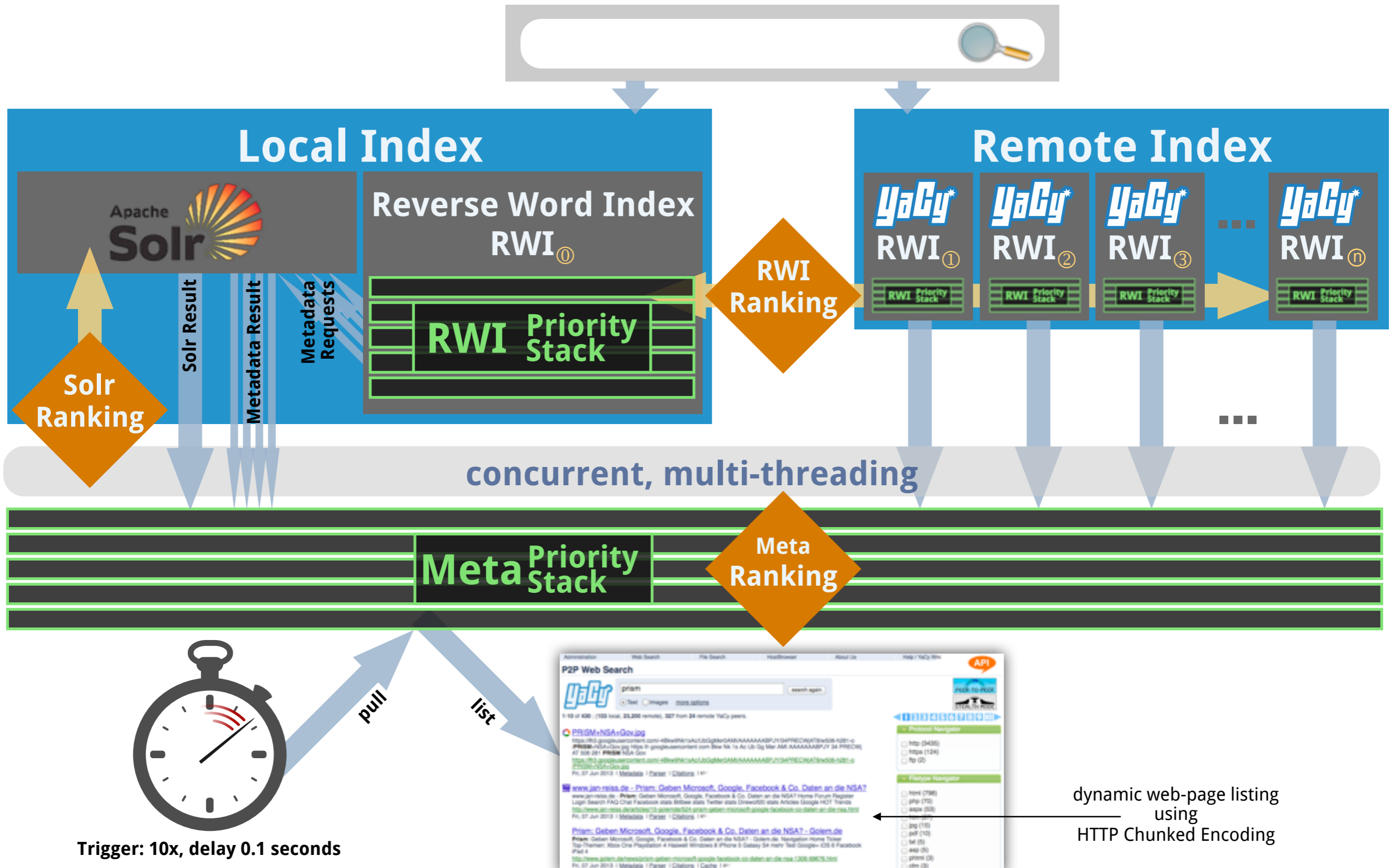
The Filter Query is attached to every query. Use this to statically add a selection criteria to reduce the set of results. Example: "http_unique_b:true AND www_unique_b:true" will filter out all results where urls appear also with/without http(s) and/or with/without 'www.' prefix. To find appropriate fields for this query, see the [YaCy Solr Schema](#). Warning: bad expressions here will cause that you don't have any search result!

fq=

Set Filter Query

Re-Set to default

YaCy verteilte Suche + Ranking

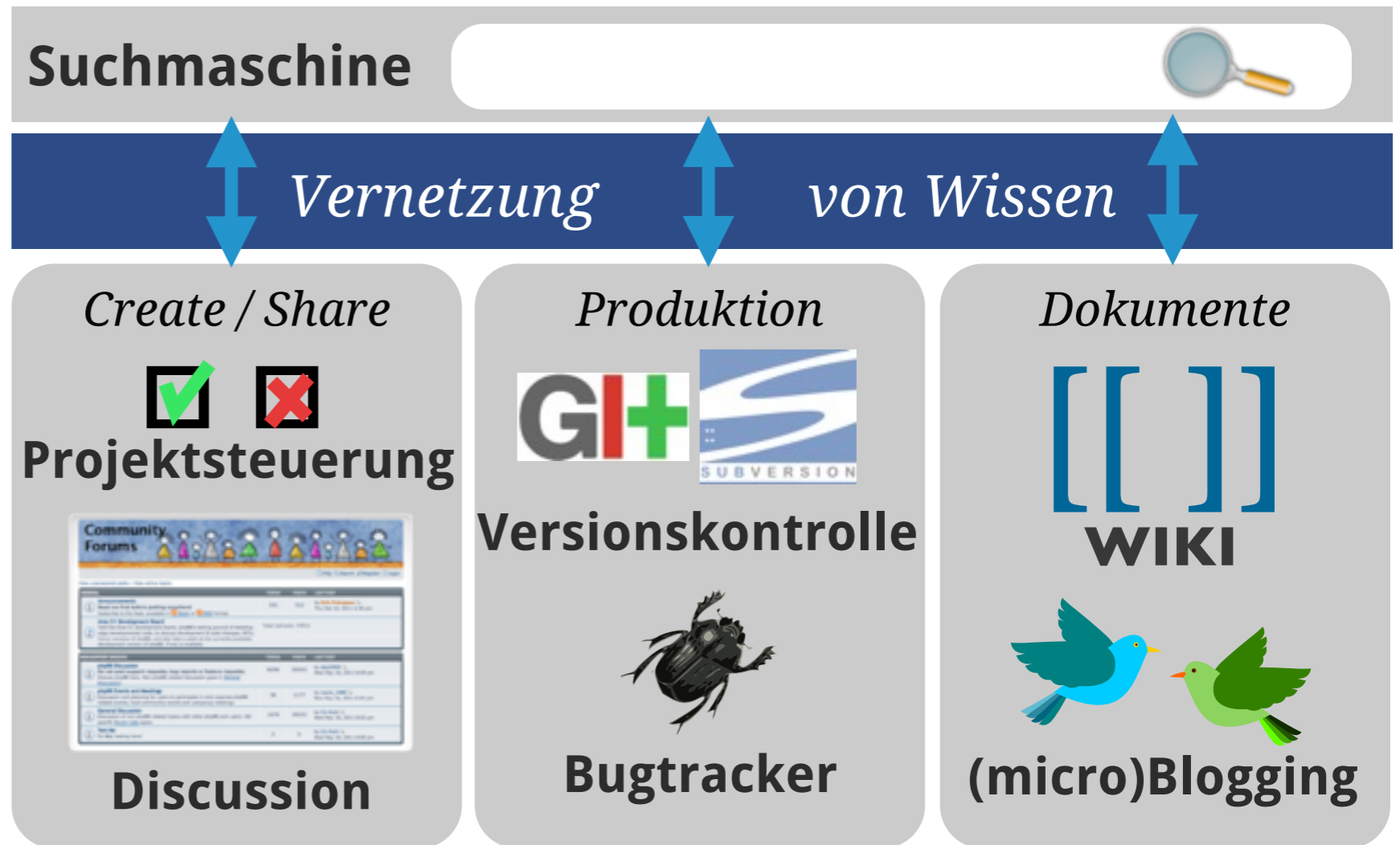


Information Retrieval mit YaCy

Michael Christen
mc@yacy.net, <http://yacy.net>



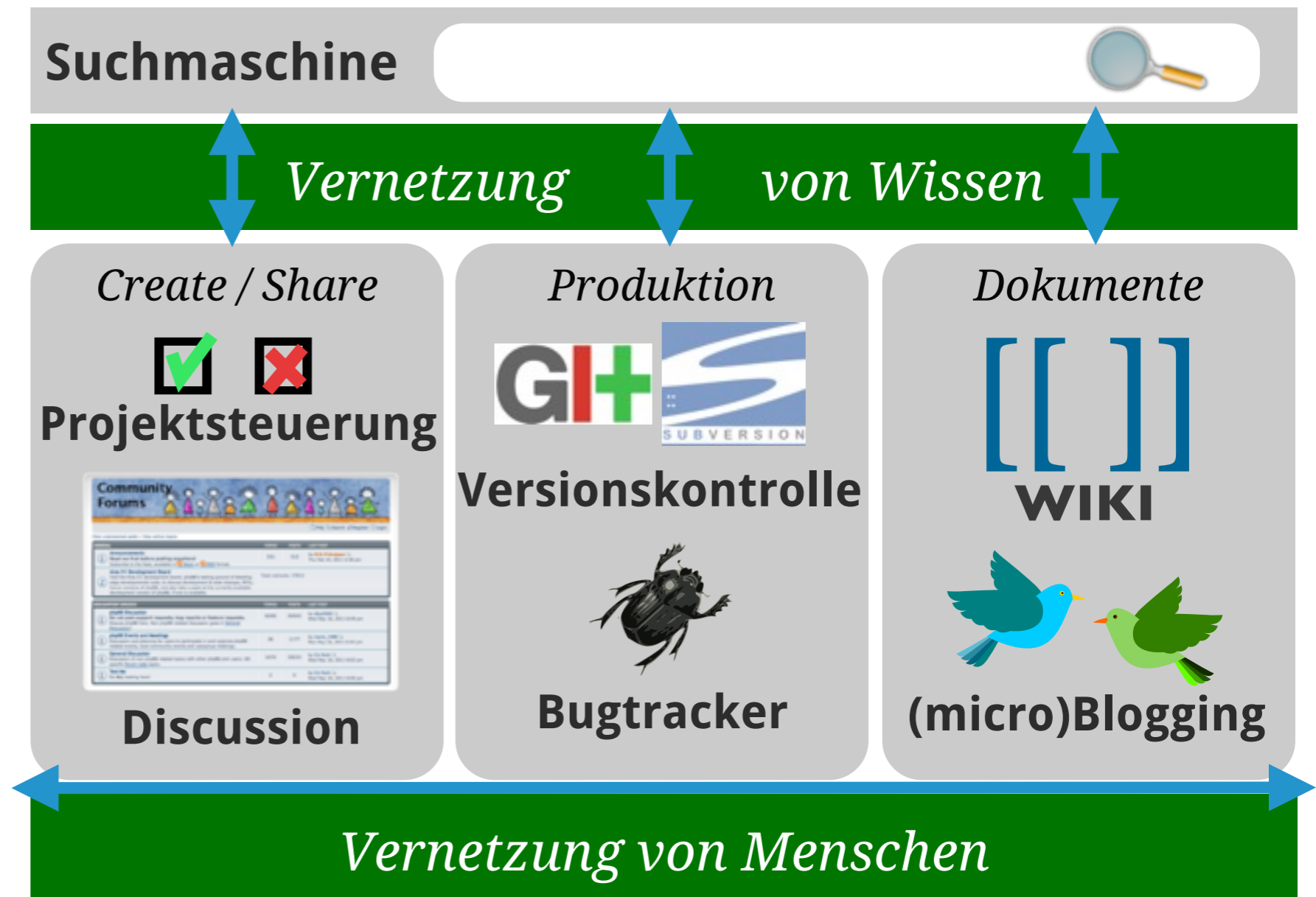
persönliche
Suchmaschine
eigene oder
fremde
Datenbestände



**persönliche
Suchmaschine**
*eigene oder
fremde
Datenbestände*

Vorteile im Unternehmen:

- Information ist unabhängig vom Ablagesystem sichtbar
- Gemeinsame Navigation unterstützt Vernetzung
- Nutzer wählen das optimale System zur Ablage



Technologische Vernetzung

„wie setze ich Technik ein um Wissen zu generieren?“

Soziotechnische Vernetzung

„wie gehen Menschen mit Technik um?“

Persönliche Suchmaschine **eigene oder fremde Daten**

**persönliche
Suchmaschine**
*eigene oder
fremde
Datenbestände*



Intranet- und Filesuche

*Konsolidierung
der Datenablage
(ftp/smb-Suche)*



Nachrichtendienst

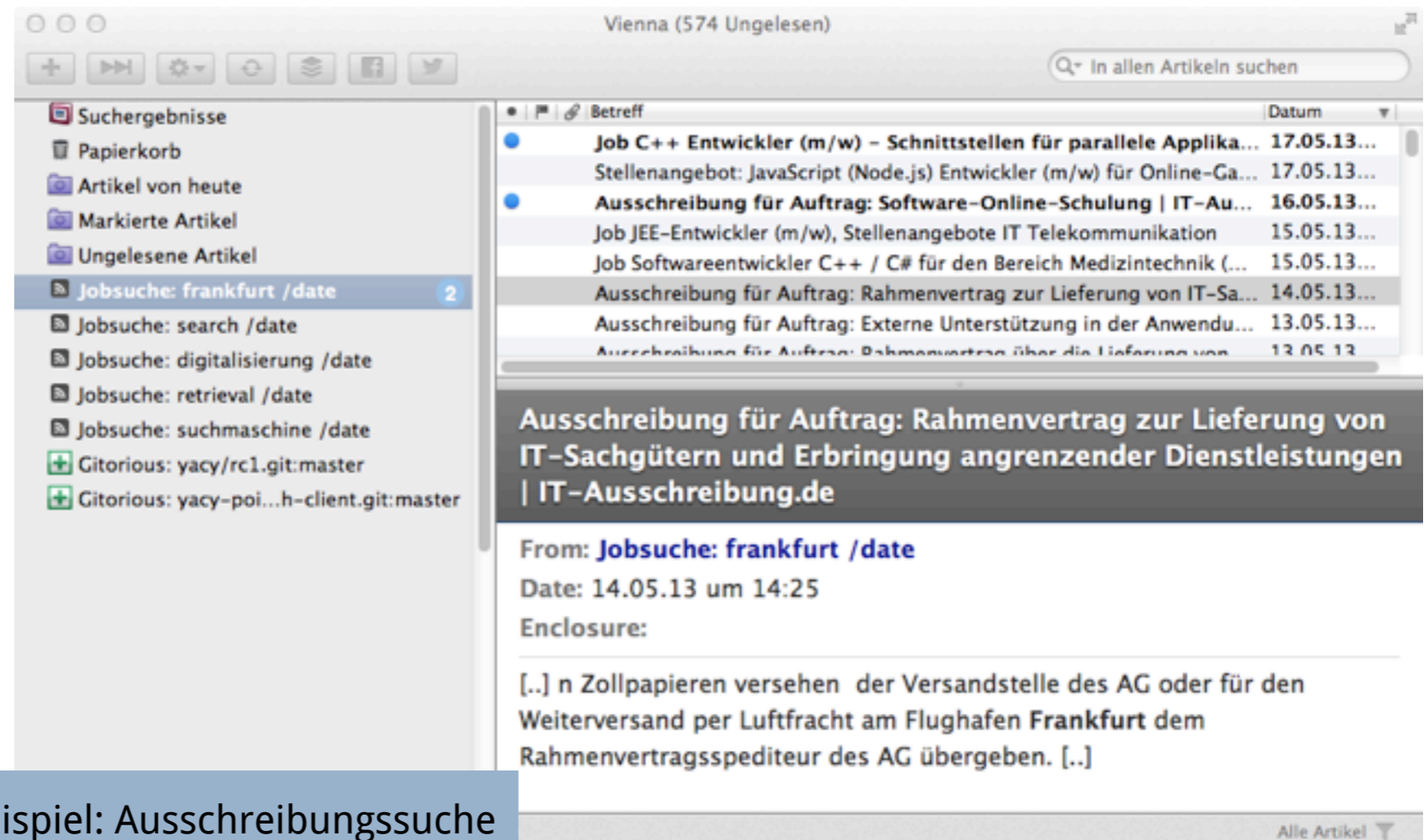
*automatisches
Suchen ohne
Suchanfrage*

Vorteile im Unternehmen:

- Zeitnahe Reaktionsfähigkeit
- Alerts für Aktivitäten der Mitarbeiter (im Intranet)
- Alerts für Aktivitäten der Konkurrenz (im Einsatz in einer Websuche)

Funktionsweise:

- Jede Suche kann ein RSS Nachrichtenstrom sein
- Suchergebnisse können nach Aktualität geordnet werden
- Suchergebnisse können automatisch weiterverarbeitet werden (RSS Reader, Alerts, u.s.w.)



Beispiel: Ausschreibungssuche

**Zitate und
Plagiate**
*Text-Überein-
stimmungen
autom. finden*

Zu jedem Suchtreffer die zitierenden Quellen finden


<http://www.welt.de/vermischtes/article116900466/Flut-hat-in-Passau-800-Haeuser-beschaedigt.html>

findet in:

<http://localhost:8090/api/citation.html?hash=6A65RsaCJDpA>

22 identische Sätze in:

http://www.focus.de/panorama/welt/tid-31632/live-ticker-zum-hochwasser-flut-welche-rolle-gen-norden-bitterfeld-evakuiert-zehntausende-menschen-fliehen_aid_1005575.html

 [Hochwasser-Ticker : Flut hat in Passau 800 Häuser besch...](#)
[Nachrichten Panorama - DIE WELT](#)
Hochwasser-Ticker
<http://www.welt.de/vermischtes/article116900466/Flut-hat-in-Passau-800-Haeuser-beschaedigt.html>
Fri, 07 Jun 2013 | [Metadata](#) | [Parse](#) | [Citations](#) | ...

**Zitate finden zu
jedem Treffer**

Document Citations for <http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html>

Similar documents from different hosts:

- http://www.focus.de/panorama/welt/tid-31632/live-ticker-zum-hochwasser-flutwelle-rolлт-gen-norden-bitterfeld-evakuiert-zehntausende-menschen-fliehen_aid_1005595.html

List of Sentences in <http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html>

0 Hochwasser-Ticker : Große Sorge um ostelbischen Teil Magdeburgs - Nachrichten Panorama - DIE WELT

http://www.focus.de/panorama/welt/tid-31632/live-ticker-zum-hochwasser-flutwelle-rolлт-gen-norden-bitterfeld-evakuiert-zehntausende-menschen-fliehen_aid_1005595.html makes 22 citations: of <http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html>

[Ministerpräsidenten Ali Larayedh brachte zu seinem Besuch bei Bundeskanzlerin Angela Merkel \(CDU\) zwei Tonnen Datteln für die Männer und Frauen in den Flutgebieten mit.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Er wolle nicht direkt von Fehlern sprechen, aber: 'Ich glaube, es waren langfristige Fehlentwicklungen, die auch teilweise vorher Experten mitgetragen haben', sagte der österreichische Wirtschaftsminister Reinhold Mitterlehner \(ÖVP\) vor dem Treffen der EU-Energieminister in Luxemburg.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Nachdem in Sachsen-Anhalt viele Menschen wegen des Hochwassers ihre Wohnungen verlassen mussten, hält die Polizei dort verstärkt Ausschau nach Plünderern.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Noch sei in keine Wohnung eingebrochen worden, sagte eine Sprecherin des Innenministeriums in Magdeburg und bestätigte damit einen Bericht der 'Magdeburger Volksstimme'.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[In den besonders stark von den Fluten betroffenen Gebieten an der Saale patrouillierten seit Tagen Polizeikräfte.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Juli\) der Kinderhilfsorganisation 'Ein Herz für Kinder' überreichen.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[In Ungarn steigt der Wasserstand der Donau weiter an.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Tausende Menschen wurden am Freitag aufgefordert, sofort ihre Häuser zu verlassen und sich in Sicherheit zu bringen.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Der Grund: Wegen Sicherungsmaßnahmen am Lober-Leine-Kanal \(Seelhausener See\) erhöht sich die Gefahr eines Wassereintruchs in den Goitzschese.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Sollte Wasser unkontrolliert in den See laufen, könnten Teile von Bitterfeld überflutet werden.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Dort stieg das Wasser in der Nacht bis auf 9,88 Meter.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Am Freitagmittag soll der Höchststand erreicht werden.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Stündlich nimmt der Wasserstand um knapp einen Zentimeter zu.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Normal sind für die Elbe in Magdeburg knapp zwei Meter.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Am Morgen ging der Pegelstand in Halle-Trotha auf 7,45 Meter zurück.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Das waren fast zehn Zentimeter weniger als in der Nacht.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Die Hochwasser-Katastrophe in Deutschland führt auch zu Behinderungen im Fernbahnverkehr.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Auch die Fahrzeiten auf der Strecke Magdeburg-Leipzig und in der Gegenrichtung verlängern sich.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Die Bahn hat eine kostenlose Hotline zu den Auswirkungen des Hochwassers auf den Bahnverkehr eingerichtet.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[An den Spendentelefonen sollen beispielsweise Kati Witt und der Moderator Florian Silbereisen sitzen.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

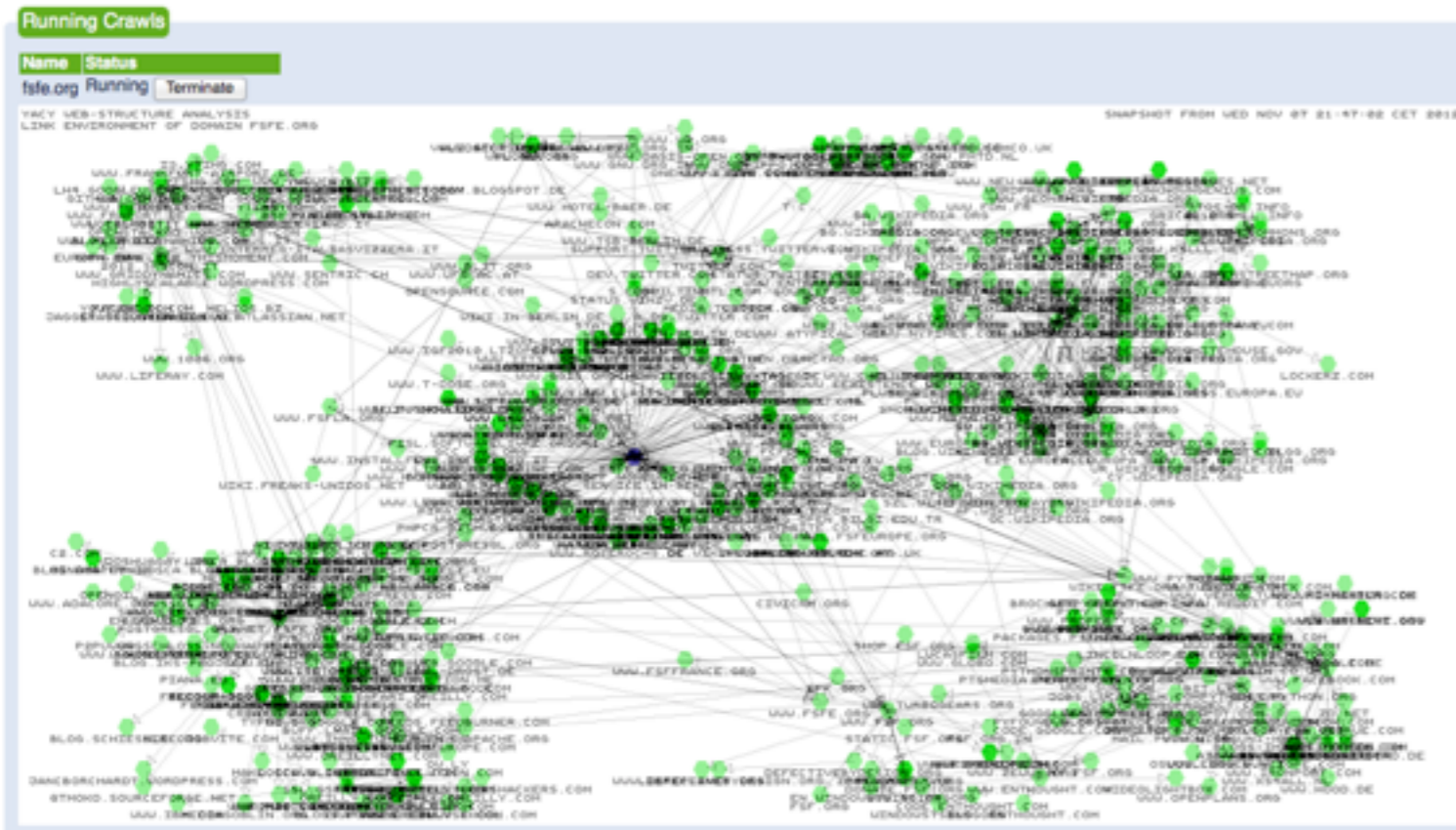
[Auch die Ministerpräsidenten der betroffenen Länder im MDR-Sendegebiet, Stanislaw Tillich \(Sachsen, CDU\), Reiner Haseloff \(Sachsen-Anhalt, CDU\) und Christine Lieberknecht \(Thüringen, CDU\), haben sich angekündigt.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

[Inka Bause und Axel Bulthaupt werden die Sendung moderieren.](http://www.welt.de/vermishtes/article116900466/Grosse-Sorge-um-ostelbischen-Teil-Magdeburgs.html)

SEO & Web-Admin Tools *Server durchstöbern, Struktur analysieren*

Funktionen:

- Die Datenstruktur fremder Server durchstöbern und Strukturen analysieren
- Tote Links aufdecken
- Visualisierung von Verlinkungsstrukturen
- Durchsuchbarkeit des eigenen Webauftritts testen



Host Browser

Browse the index of 4,105 documents. Enter a host or an URL for a file list or select one of a [list of hosts](#).

Host/URL: Browse Host

Host List

fsfe.org	637/569/1 URLs	en.wikipedia.org	344/268/1 URLs
twitter.com	250/29/13 URLs	blogs.fsfe.org	204/177 URLs
identi.ca	139/131/2 URLs	www.computerworld.com.au	134/131 URLs
sankleif.wordpress.com	121/116 URLs	grical.org	118/66 URLs
computerfloss.com	113/109 URLs	blog.padowi.se	106/103/1 URLs
www.fsfla.org	104/101/2 URLs	www.suedtirolerland.it	103/98/2 URLs
typo3.org	103/100/1 URLs	wiki.fsfe.org	104/96 URLs
www.cabinetoffice.gov.uk	100/89/3 URLs	www.fsf.org	97/90/1 URLs
archive.org	94 URLs	www.fdn.fr	84/60/1 URLs
de.wikipedia.org	78/52/2 URLs	leana.de	78/75 URLs
openoil.net	72/69 URLs	honk.sigxcpu.org	70/64/1 URLs
seravo.fi	68/63 URLs	hircus.wordpress.com	67/63 URLs
blog.ks-project.eu	65/62 URLs	directory.fsf.org	65/60 URLs
www.gnu.org	63/56 URLs	www.mediamatic.net	59/58 URLs
news.swpat.org	57/55 URLs	lhuffpost.com	54 URLs
www.youtube.com	51/43 URLs	losca.blogspot.fi	47/44/1 URLs
www.techcast.com	47/47 URLs	www.linuxtag.org	45/35/2 URLs
www.installfest.info	45/42/1 URLs	www.adacore.com	44/40/1 URLs
www.huffingtonpost.com	45 URLs	softwarefreedomday.org	43/40/1 URLs
www.tis.bz.it	43/38 URLs	www.gag.com	42/36 URLs
ec.europa.eu	39/36/2 URLs	florian.wordpress.com	40/36 URLs
2010.mil.info	39/38 URLs	opensource.com	37/34 URLs
micro.systemsavivour.com	35/33 URLs	www.dartlang.org	35/23 URLs
www.welt.de	35 URLs	www.4shared.net	35 URLs
www.lockergnome.com	34/31 URLs	www.kiberpipa.org	33/25/1 URLs
summit.ubuntu.com	32/29 URLs	linuxwochen.at	32/28 URLs
news.yahoo.com	32 URLs	www.flossk.org	31/28 URLs
www.gnu.org.in	31/26 URLs	www.linuxpromagazine.com	28/25 URLs

Web Analytics
*Webseiten an die
Anforderungen
der User
anpassen*



- ➔ erkennen, was die User von der Website erwarten
- ➔ Keywordanalyse

http://localhost:8090/api/timeline_p.xml?from=20130101000000&to=20131231000000&data=queries&head=30&period=1y
-> Top-100 der meistgesuchten Begriffe auf fsfe.org im Jahr 2013

```
<event time="20130101000000" isPeriod="true" duration="3153600000" count="283" type="query">"free symbian apps"</event>  
<event time="20130101000000" isPeriod="true" duration="3153600000" count="118" type="query">cap</event>  
<event time="20130101000000" isPeriod="true" duration="3153600000" count="90" type="query">android</event>  
<event time="20130101000000" isPeriod="true" duration="3153600000" count="70" type="query">test</event>  
<event time="20130101000000" isPeriod="true" duration="3153600000" count="65" type="query">l</event>  
<event time="20130101000000" isPeriod="true" duration="3153600000" count="64" type="query">pdf</event>  
<event time="20130101000000" isPeriod="true" duration="3153600000" count="51" type="query">yacy</event>  
<event time="20130101000000" isPeriod="true" duration="3153600000" count="48" type="query">teen</event>  
<event time="20130101000000" isPeriod="true" duration="3153600000" count="45" type="query">root</event>  
<event time="20130101000000" isPeriod="true" duration="3153600000" count="31" type="query">http://www.generic4you.com/17227-viagra</event>  
<event time="20130101000000" isPeriod="true" duration="3153600000" count="28" type="query">sex</event>  
<event time="20130101000000" isPeriod="true" duration="3153600000" count="28" type="query">gnu</event>
```




Vielen Dank fürs Zuhören

Dipl. Inf. Michael Christen

✉ mc@yacy.net

🐦 @0rb1t3r

Links

YaCy Home Page

<http://yacy.net>

(Downloads, YaCy Forum + YaCy Wiki + Bugtracker dort verlinkt)

Source Code GIT

<https://gitorious.org/yacy/rc1/>

Tutorial Videos

<https://www.youtube.com/user/YaCyTutorials/>

Social Media

https://twitter.com/yacy_search

<https://www.facebook.com/yacy.search.engine>

