

10 years of „Bielefeld Academic Search Engine“ (BASE)

Dirk Pieper - Friedrich Summann
Bielefeld University Library

10 years of „Bielefeld Academic Search Engine“ (BASE) (www.base-search.net)

**Looking at the past and future of the world wide
repository landscape from a service providers
perspective**

Overview

- **Retrospect (2001 – 2014)**
- Current Situation (2015)
- Future prospective (2015 - ...)

Some Milestones (I)

- 2001 Starting point as a search engine follow-up for a metasearch system
- 2004 Official Start (FAST Data Search)
- 2005 search history, sorting
- 2006 starting participation in EU projects (DRIVER), BASE-APIs
- 2007 Introducing multilingual search
- 2008 Index > 10 million records

Some Milestones (II)

- 2009-2011 Automatic classification of OAI-DC-Metadata (funded by German Research Foundation)
- 2011 Switch to open source (Lucene/Solr, VuFind)
- 2011 Index > 30 million records, Iphone-App
- 2012 Web-App, OAI-Interface, data delivery of subject sections
- 2014 OA-boosting

BASE Digital Collections**Demonstr****BASE includes the following content sources****Basic Sea**

Deutsche \

● free con

[Mathematical Preprint Server](#), Bielefeld University, Faculty for Mathematics
(complete)

[E-journal Documenta Mathematica](#), Bielefeld University, Faculty for Mathematics
(complete)

[Online Publications](#), Bielefeld University (complete)

[Journals of the German Enlightenment](#), Bielefeld University Library (complete)

[Library Catalogue](#), Bielefeld University Library (partial)

[Electronic Dissertations](#), Bochum University (complete)

[Project Euclid](#), Cornell University (partial)

[Historical Math Monographs](#), Cornell University Library (partial)

[Internet Library of Early Journals](#), Host: Oxford University Library Services
(complete)

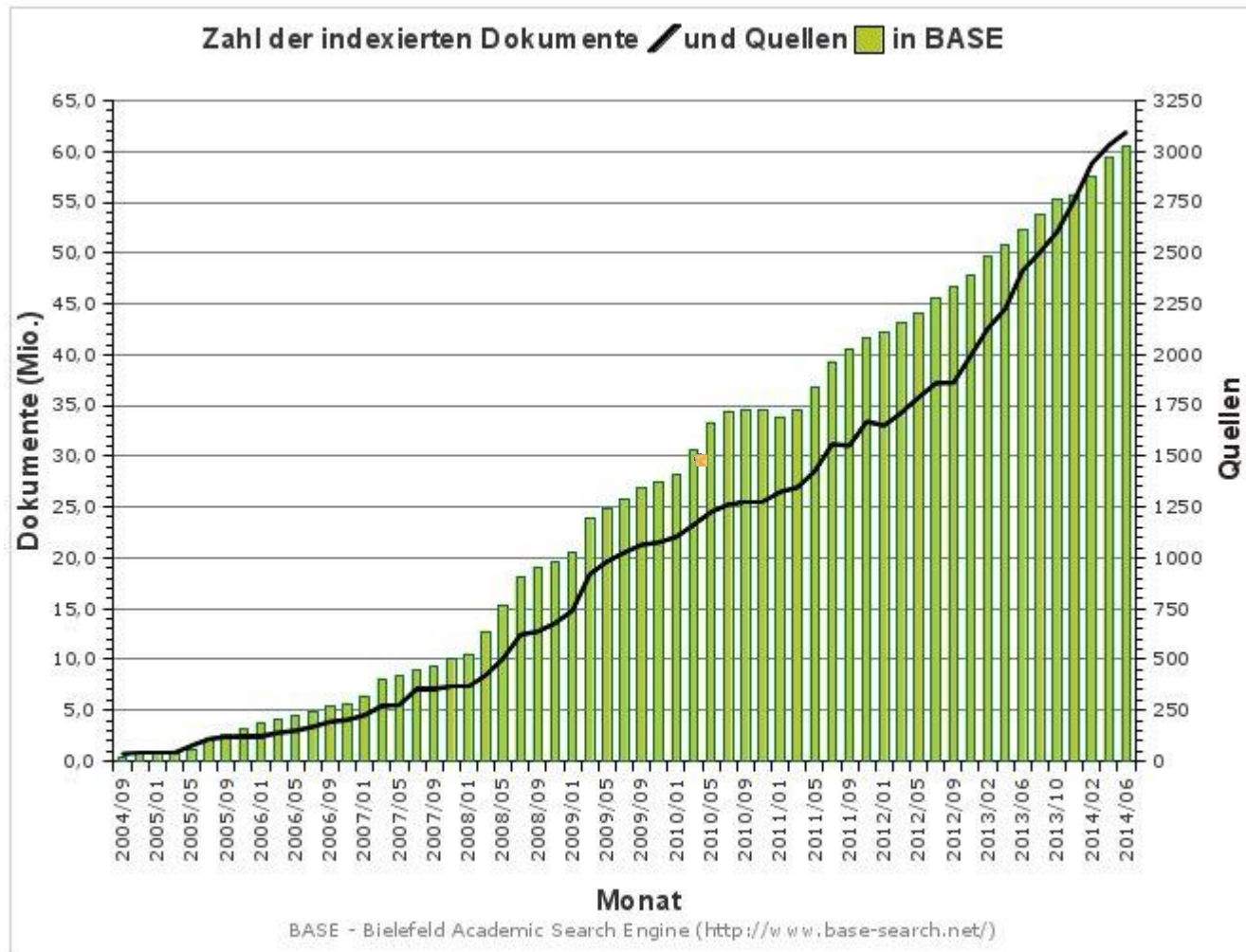
[Mathematical Collection](#), Springer Publishing House Heidelberg, Germany (partial)

[Digital Collection Mathematica](#), SUB Goettingen/GDZ (complete)

[Research Reports, BMBF](#), TIB/UB Hannover (partial)

[Historical Mathematics Collection](#), The University of Michigan (partial)

[Zentralblatt MATH](#), Host: FIZ Karlsruhe / Zentralblatt für Mathematik (partial)



Tabelle

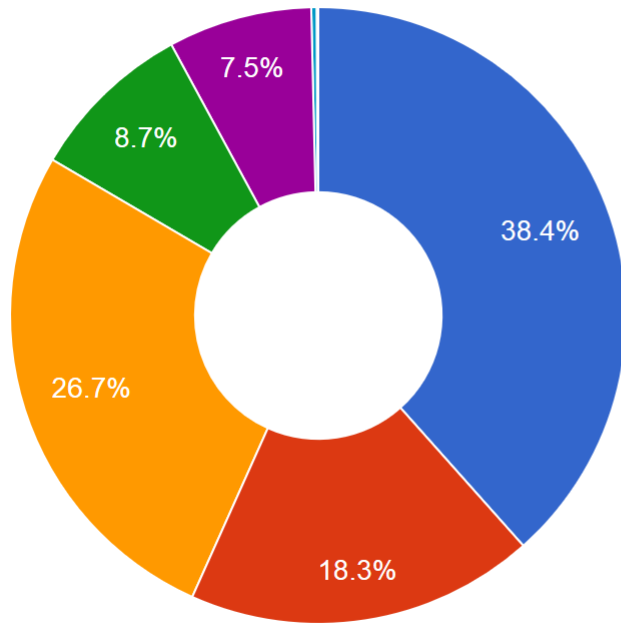
Overview

- Retrospect (2001 – 2014)
- **Current Situation (2015)**
- Future prospective (2015 - ...)

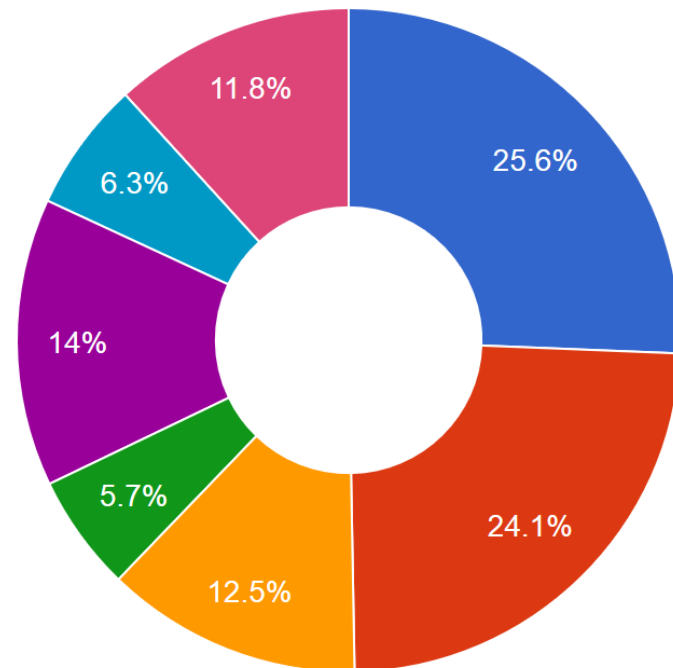
The BASE scope

- OA Repositories world-wide (institutional and subject reps)
- Academic Valuable Content
- Electronic Journals
- Aggregators (RePEc, Virtual Libraries, etc.)
- Digital Collections
- Dataset Repositories

Repository Types covered in BASE



- Institutional Repositories
- Publication Server
- Electronic Journals
- Thesis/Dissertation Servers
- Digital Collections
- Research Data
- Sonstiges



Contents:

Global Scientific Metadata harvested from
Repositories

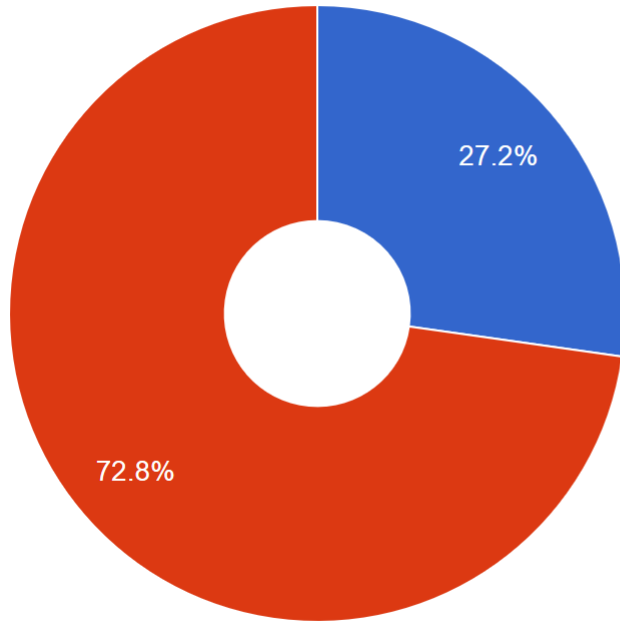
Technology:

- Index (Solr)
- Search Interface (VuFind)
- APIs (HTTP, OAI)

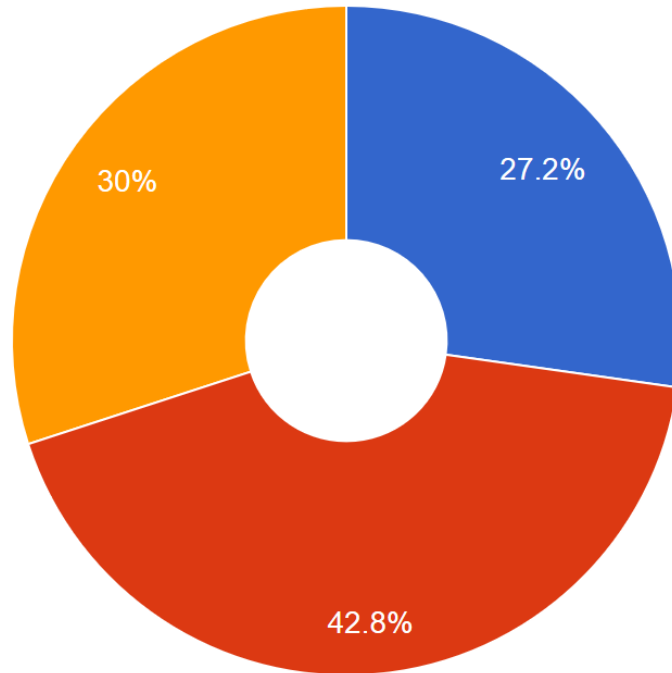
BASE harvests, aggregates, enriches and exposes
OAI Metadata



Open Access Status



- Open Access (definitive)
- unknown



**BASIC
SEARCH**

ADVANCED
SEARCH

HELP

BROWSING

SEARCH
HISTORY

OVER 70 MILLIONS DOCUMENTS

Your search

repositories

Entire Document

- Boost open access documents
- Verbatim search
- Additional word forms
- Multilingual synonyms (Eurovoc Thesaurus)

Find

Currently in BASE: 73,740,686 Documents of 3,531 [Content Sources](#)

[About BASE](#)    | [Contact](#) | [BASE Lab](#) | [Imprint](#)

© 2004-2015 by [Bielefeld University Library](#)

Search powered by [Solr](#) & [VuFind](#).

 [Suggest Repository](#)
 [BASE Interfaces](#)

Hit List

1. HKUST Institutional Repository

Title: HKUST Institutional Repository

Author: Wang, Gabriela K. W.

Description:

Hit List

1. Wikibooks: Git/Repository on a USB stick

 Open Access

Title: Wikibooks: Git/Repository on a USB stick

Description: Instead of having to resort to a hosting company to store your central repository or to rely on a central server or internet connection to contribute changes to a project it s quite possible to use removeable memory to exchange and update local **repositories**. The basic steps are Mount the removable memory on a pre determined path Setup a bare re... [+ Show all](#)

Document Type: Book

Language: eng

URL: http://en.wikibooks.org/wiki/Git/Repository_on_a_USB_stick

Content Provider: **WikiBooks - Open-content textbooks**

2. Rep

[Detail View](#) [Email this](#) [Export Record](#) [Add to Favorites](#) [Check in Google Scholar](#)

Title:

Author:

Publ

Year

2. RepositoriUM. Institutional Repository. Minho University, Braga, Portugal: A university repository where a mandate to deposit, financial incentives and strong advocacy can...

 Open Access

Title: RepositoriUM. Institutional Repository. Minho University, Braga, Portugal: A university repository where a mandate to deposit, financial incentives and strong advocacy can transform an Institutional Repository's population

Author: **Proudman, V.M.**

Facts about BASE

- 3531 Repositories included
- From 102 Countries world-wide
- Ca. 73 Mill. Documents/Objects
- Ca. 70 % Open Accessible
- Ca. 10.3 Mill. Documents enriched with DDC-Codes (Dewey)

Harvesting Environment (Based on OAI-PMH)

- 5979 Repositories harvested
 - 4832 active
 - 3531 indexed
 - 1147 deprecated
- 184 Mill. Records
 - 111 Mill. unique
 - 73 Mill. indexed
- 1.03 Terabyte of Data
- 2712 Cronjobs (weekly)

Khabsa M, Giles CL (2014) The Number of Scholarly Documents on the Public Web.
PLOS ONE 9 (5) DOI: 10.1371/journal.pone.0093949



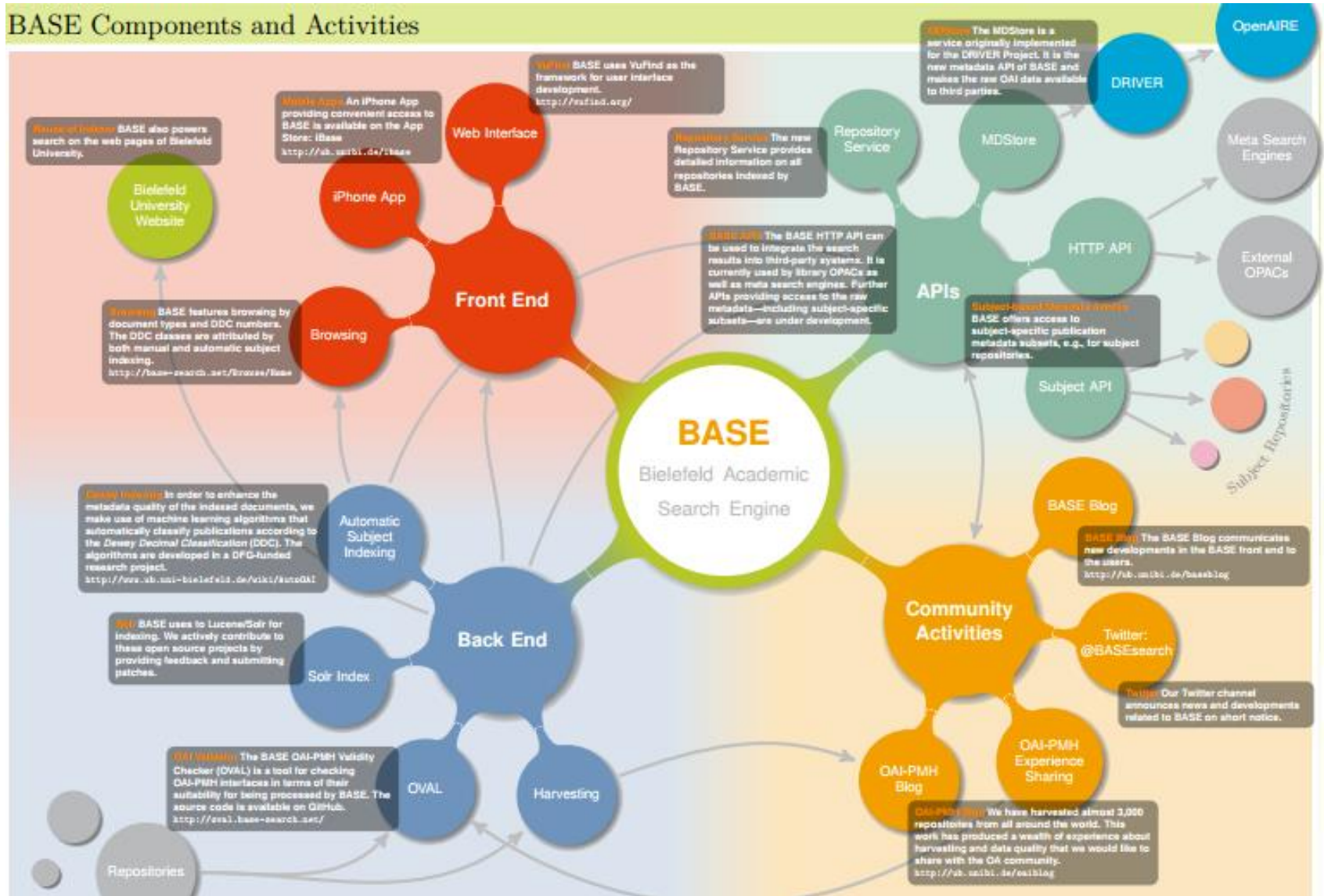
Figure 2. Relative number of documents by scholarly search engines and databases. Total and Google Scholar are estimates.
doi:10.1371/journal.pone.0093949.g002

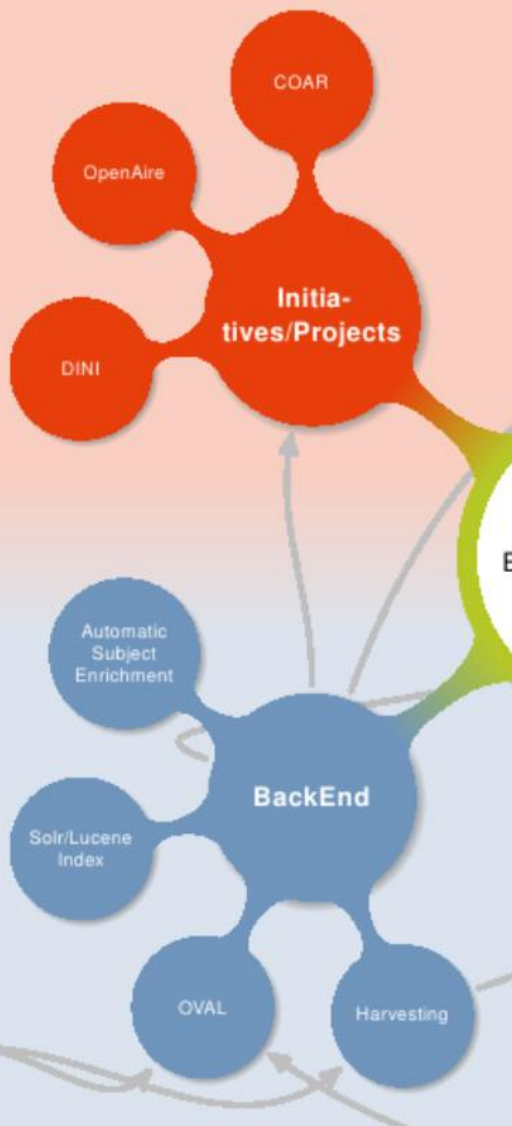
Basic Principles



- **B**ibliographic Expertise
- **A**utomatic methods
- **S**imple Infrastructure
- **E**fficient Strategies

BASE Components and Activities





迈向下一代 OAI 服务提供商，以 BASE 为例

Dirk Pflaeg, Friedrich Summann, Bernd Fehling, Renata Mirenga, Sebastian Wolf, and Matthias Lüscher
www.base-search.net

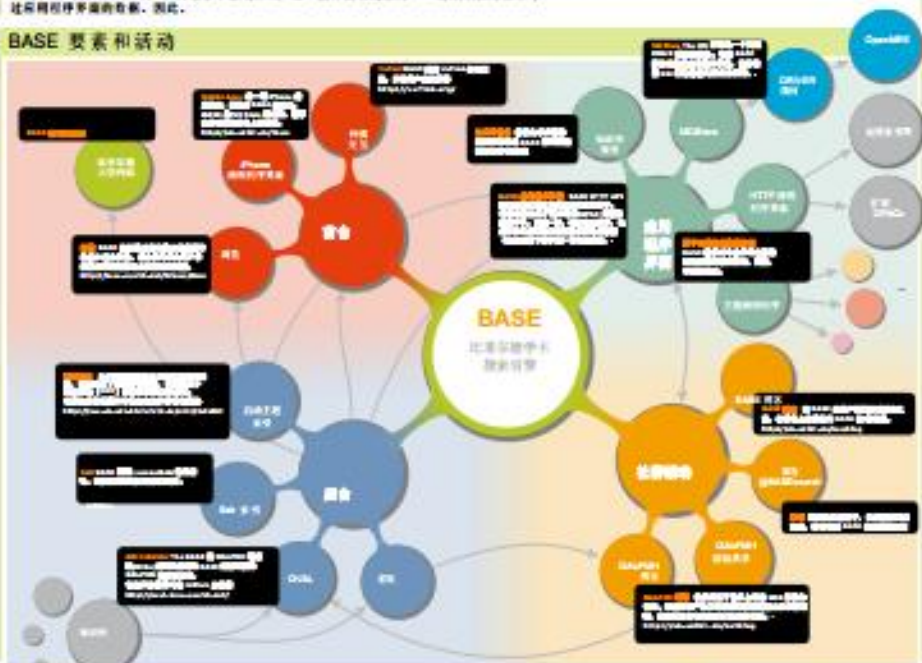
简介

OAI 服务提供商已经成了学术交流的良好帮手。不过，学术网络的发现并不成熟。已有许多提供商努力完成数据集成并且提供统一交互界面使搜索便利。然而它们一成不变。服务提供商可以挖掘加以利用支持的加强服务的潜力。像丰富集成、个性化、关联搜索等。它给服务提供商集成数据新提供者—提供集成后的—以及整合中强化后的一级过应用程序界面的有数。因此，

所以提供商在知识库中进行的数据检索。以及深化服务的不断实现。通过 BASE 的几项活动，有比率为 OAI 服务提供商建设，以满足这些需求。自 2004 年起，比萨大学图书馆开始地作为为计划的服务提供商 BASE 以来，开发了一系列有关集成、个性化、定向 OAI-PMH 索引数据。并提供世界上最完整的 OAI 连接数据的集合等。

BASE 为其他服务提供商，分享数据接口和索引等等。还有，我们开发 1 德国科学基金会资助的 1 基于学习的文本自动分类系统。应用 OAI 的索引和元数据记录进行主题索引自动提取。能够提供主题下的特字主题的集成访问功能。...

BASE 要素和结构



BASE 大事记



- 新增超过 10 个图书馆合作
- 新增超过 10 个数据集
- 新增七十多个新的地址
- 新增中文及日文语言内容数据集
- 新增多语言、多格式数据集

BASE 工作组

Dirk Pflaeg: 项目合作者 / 技术专家
Friedrich Summann: 比萨大学图书馆 IT 部门主管
Bernd Fehling: 海防库成员
Renata Mirenga: 青年研究员
Sebastian Wolf: 电子服务部门
Matthias Lüscher: 程序员 / 系统管理员、文本数据



www.base-search.net

BASE restrictions for including repositories

- OAI-PMH interface (available)
- OA is supported (at least some documents should be open accessible world-wide!)
- Valid DC Metadata (basic set of field contents)
- Operable URL for the document landing page

Main Issues to Avoid

- Wrong Document URLs
- Empty Records
- Invalid XML delivered
- Crashing Harvesting Processes
- Resumption Token Handling Fails
- No Incremental Delivery
- Changing OAI-PMH or homepage URLs without dissemination/redirect
- No OAI-PMH Interface Configuration

Main Issues Preferred

- OA Status delivered (on repository or document level)
- Metadata Guidelines compatible
 - Vocabularies used (for type, language, date, classification etc.)
- (Parallel) English end-user interface
- Citation/Abstract information delivered
- Available repository contact information
- Visible in Registries (OpenDOAR, openarchives ...)

The **10** Biggest Misunderstandings around OAI-PMH

- OAI-PMH means ‚Everything is Open Access‘
- Persistent Identifiers are persistent
- Link to the Document page is not necessary
- Configuration is not needed
- Checking the Service is needless
- DublinCore is simple but sufficient
- OpenAccess Status Information is needless
- End-User Interface is not necessary
- Personal Email Adress is not needed
- Incremental Harvesting is sufficient

- Dear Sys Admin,

we have harvested and indexed the OAI-PMH interface of your repository since longer.

But when checking 1 problem:

The oai_dc records (**Mail-Prepper** for example



- Std
- No records match
- Tomcat-Error
- Error ResumptionToken Handling
- Error NoOA Documents
- Error Empty Responses
- Handle not registered (DSpace)
- Handle not configured (DSpace)

Dear Friedrich,
Thank you very much
me if it is ok for you.
an absolute one.
Regards
Stéphanie

BASE Usage (May 2015)

Länder		Seiten	Zugriffe	Bytes	
	Germany	de	154,631	665,922	19.12 GB
	China	cn	80,942	193,681	6.23 GB
	United States	us	78,107	142,989	7.40 GB
	France	fr	56,941	168,232	4.79 GB
	Canada	ca	42,242	59,315	2.16 GB
	Great Britain	gb	11,401	36,616	1.51 GB
	Poland	pl	9,618	42,138	1.39 GB
	Netherlands	nl	9,378	17,740	761.28 MB
	Switzerland				
	Austria				
	Peru				
	Spain				
	Belgium				
	Mexico				
	Ukraine				
	India				
	Colombia				
	Brazil				
	Italy				
	Argentina	ar	4,457	20,568	558.23 MB

Monat	Unterschiedliche Besucher	Anzahl der Besuche	Seiten	Zugriffe	Bytes
Jan 2015	39,030	75,098	670,532	2,108,535	67.32 GB
Feb 2015	37,677	72,774	645,761	1,976,945	63.28 GB
März 2015	44,295	86,937	805,643	2,375,123	81.41 GB
Apr 2015	44,017	84,034	882,180	2,416,771	86.39 GB
Mai 2015	42,028	81,204	584,379	1,961,255	59.67 GB
Juni 2015	290	308	1,368	4,624	123.26 MB
Juli 2015	0	0	0	0	0
Aug 2015	0	0	0	0	0
Sep 2015	0	0	0	0	0
Okt 2015	0	0	0	0	0
Nov 2015	0	0	0	0	0
Dez 2015	0	0	0	0	0
Total	207,337	400,355	3,589,863	10,843,253	358.18 GB



Overview

- Retrospect (2001 – 2014)
- Current Situation (2015)
- **Future prospective (2015 - ...)**

Current developments

- Analysis of alternative (OAI-PMH) metadata formats
- Prototypical expansion of metadata workflows (use case: SCOAP3-repository)
- Increasing the OA percentage within BASE by detailed analysing of metadata (setSpecs, rights-information)
- Working on grant proposal with Helmholtz Association and German National Library for an ORCID Claiming Service
- Preparing resolving service for several identifiers

BASE Future Strategies

- Optimizing the Technical Platform
- Link Resolver Service
- Data Enrichment (Linked Open Data Strategies)
- Big Data Activities

Optimizing the Technical Platform

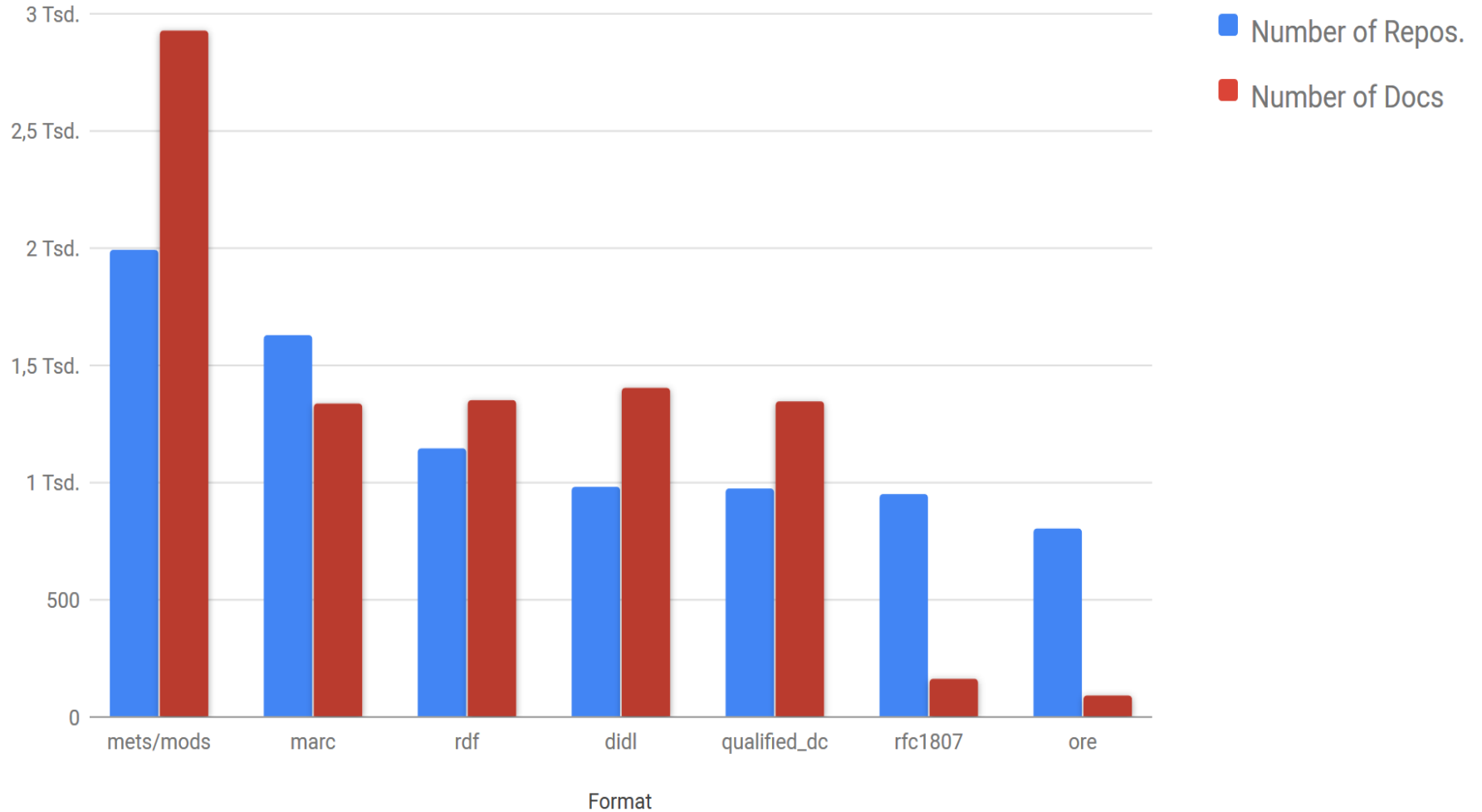
- Optimizing the VM Server Structure
- Distribution of Applications
- Optimizing the Lucene/Solr-Search Engine Environment
- Multi-node System

Data Enrichment (Linked Open Data Strategies)

Fundamental Aspect:

- More **Detailed Metadata Formats** have the potential to provide more detailed information
- But: It depends on the background quality


Metadata Formats in OAI-PMH Repositories



Link Resolver Service

- OAI Identifier
- DOI/Handle/ISSN/ISBN/URN/PMCID
- Author/Organization/Funder IDs
- Bibliographic Metadata

• Issue: Rights and Licences Normalization

SUBGoettinger	 DE	ceu	181	Scan	<p>5x http://creativecommons.org/licenses/by-sa/3.0</p> <p>5x http://creativecommons.org/licenses/by-nc/2.5/</p> <p>5x http://creativecommons.org/licenses/by-nc/2.0/</p> <p>4x http://goedoc.uni-goettingen.de</p> <p>4x http://info:eu-repo/semantics/http://creativecommons.org/licenses/by-nd/3.0/de/</p> <p>4x https://creativecommons.org/licenses/by/4.0</p> <p>4x http://creativecommons.org/licenses/by/2.0/uk/legalcode</p> <p>4x http://creativecommons.org/licenses/by-nc-nd/3.0/deed.de</p> <p>4x http://creativecommons.org/licenses/by-nc/3.0/us/</p> <p>3x http://iopscience.iop.org/1367-2630/14/10/103012/</p> <p>3x http://creativecommons.org/licenses/by/2.5/http://creativecommons.org/licenses/by/2.5/</p> <p>3x 719460905</p> <p>3x 9</p> <p>3x info:eu-repo/semantics/openAccess</p> <p>3x http://creativecommons.org/licenses/by-nc-nd/2.5/ch/deed.en</p> <p>3x http://creativecommons.org/licenses/by/2.0</p> <p>3x http://creativecommons.org/licenses/by-nc-sa/3.0</p> <p>3x http://creativecommons.org/licenses/by-nc/3.0/deed.en_US</p> <p>3x creativecommons.org/licenses/by-nc-sa/3.0/</p> <p>3x http://info:eu-repo/semantics/http://creativecommons.org/licenses/by/2.5/</p> <p>3x http://creativecommons.org/licenses/by-sa/3.0/de/deed.en</p> <p>3x http://creativecommons.org/licenses/by-sa/3.0/deed.de</p> <p>2x https://creativecommons.org/licenses/by-nc-nd/3.0/</p> <p>2x http://creativecommons.org/licenses/by-nc/2.0/de/deed.en</p> <p>2x http://creativecommons.org/licenses/by-nc-nd/3.0/at/</p> <p>2x http:// http://creativecommons.org/licenses/by/2.0</p> <p>2x 504430</p>
---------------	-----------------------------------------------------------------------------------------	-----	-----	------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

Big Data Activities

- Data Enrichment
(Linked Open Data Strategies)
- Automatic Classification
- De-Duplication/Version Detection
- Fulltext Indexing

Thank you for your attention!

- friedrich.summann@uni-bielefeld.de
- dirk.pieper@uni-bielefeld.de