

Apache Solr

An experience report

2013-10-23 - Corsin Decurtins





Apache Solr

Notes

Full-Text Search Engine
Apache Lucene Project
based on Apache Lucene

Fast
Proven and Well-Known Technology
Java based
Open APIs
Customizable
Clustering Features



Setting the **Scene**

[Search Tipps](#)

Everything

Wiki

Files

Emails

Issues

People

Companies

Projects

Resources

Books & Magazines

Include Archive

All Types

Plaza Search

Plaza Search is an integrated search engine for the Netcetera intranet. You can use Plaza Search to look for resources in the **Plaza Wiki**, **our file system**, **JIRA**, **Mailstore** and in Infostore (**people**, **projects**, **companies** and **books and magazines**).

Just put your search terms into the corresponding field at the top and use the filters on the left side to further narrow down the results.

If you want to find out what Plaza Search can do for you, make sure to check out the **Plaza Search User Documentation**.

And if something should not work as you expect it or you have an idea about how to make Plaza Search even better, make sure to let us know - either by **creating a JIRA issue** or by **sending us an email**.

Plaza Search

Notes

Full-Text Search Engine for the Intranet of Netcetera

Integrates Various Data Sources

Needs to be fast

Ranking is crucial

Simple Searching

Relevant Filtering Options

Desktop, Tables and Phones

Warum Intranet-Suchmaschinen unbrauchbar sind ...und was dagegen getan werden kann

2013-07-03 – Corsin Decurtins



[Search Tipps](#)

Everything

Wiki

Files

Emails

Issues

People

Companies

Projects

Resources

Books & Magazines

Include Archive

All Types

Plaza Search

Plaza Search is an integrated search engine for the Netcetera intranet. You can use Plaza Search to look for resources in the **Plaza Wiki**, **our file system**, **JIRA**, **Mailstore** and in Infostore (**people**, **projects**, **companies** and **books and magazines**).

Just put your search terms into the corresponding field at the top and use the filters on the left side to further narrow down the results.

If you want to find out what Plaza Search can do for you, make sure to check out the **Plaza Search User Documentation**.

And if something should not work as you expect it or you have an idea about how to make Plaza Search even better, make sure to let us know - either by **creating a JIRA issue** or by **sending us an email**.

Some Numbers

Live since

05/2012

Data since

1996

~ **275** Users

~ **500 – 2'000**

Searches per day

~ **3'000'000**

Documents

~ **40** Releases

Index Size

~ **75** GByte

Some Numbers

Notes

Very small load (only a few hundred requests per day)

The indexer agents actually produce a lot more load than the actual end users

Medium size index (at least I think)

Not that many objects, but relatively big documents

Load performance is not a big topic for us

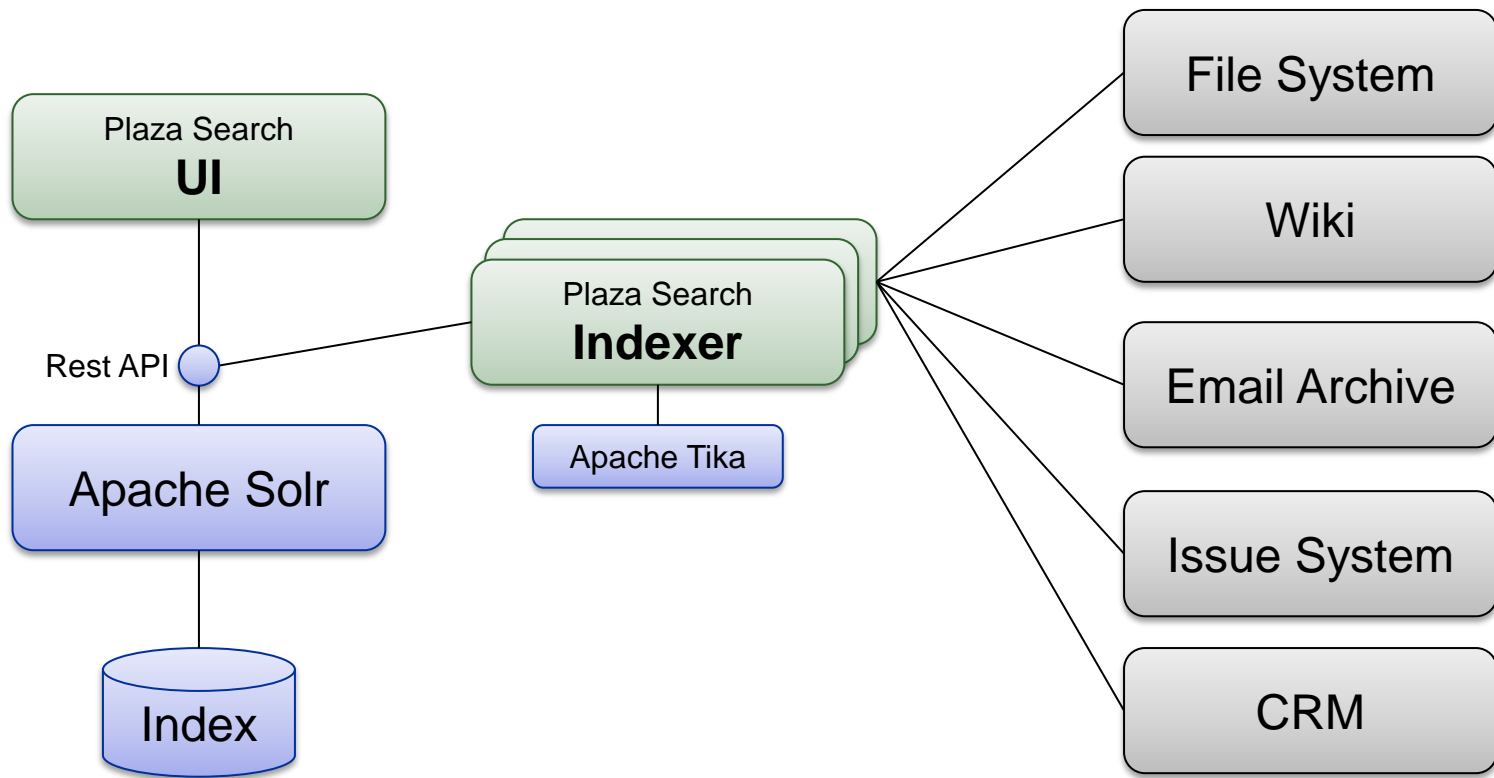
When we talk about performance, we actually usually mean response time

For us

Performance

means

Response Time



Architecture

Notes

Based on Apache Solr (and other components)

Apache Solr takes care of the text-search aspect

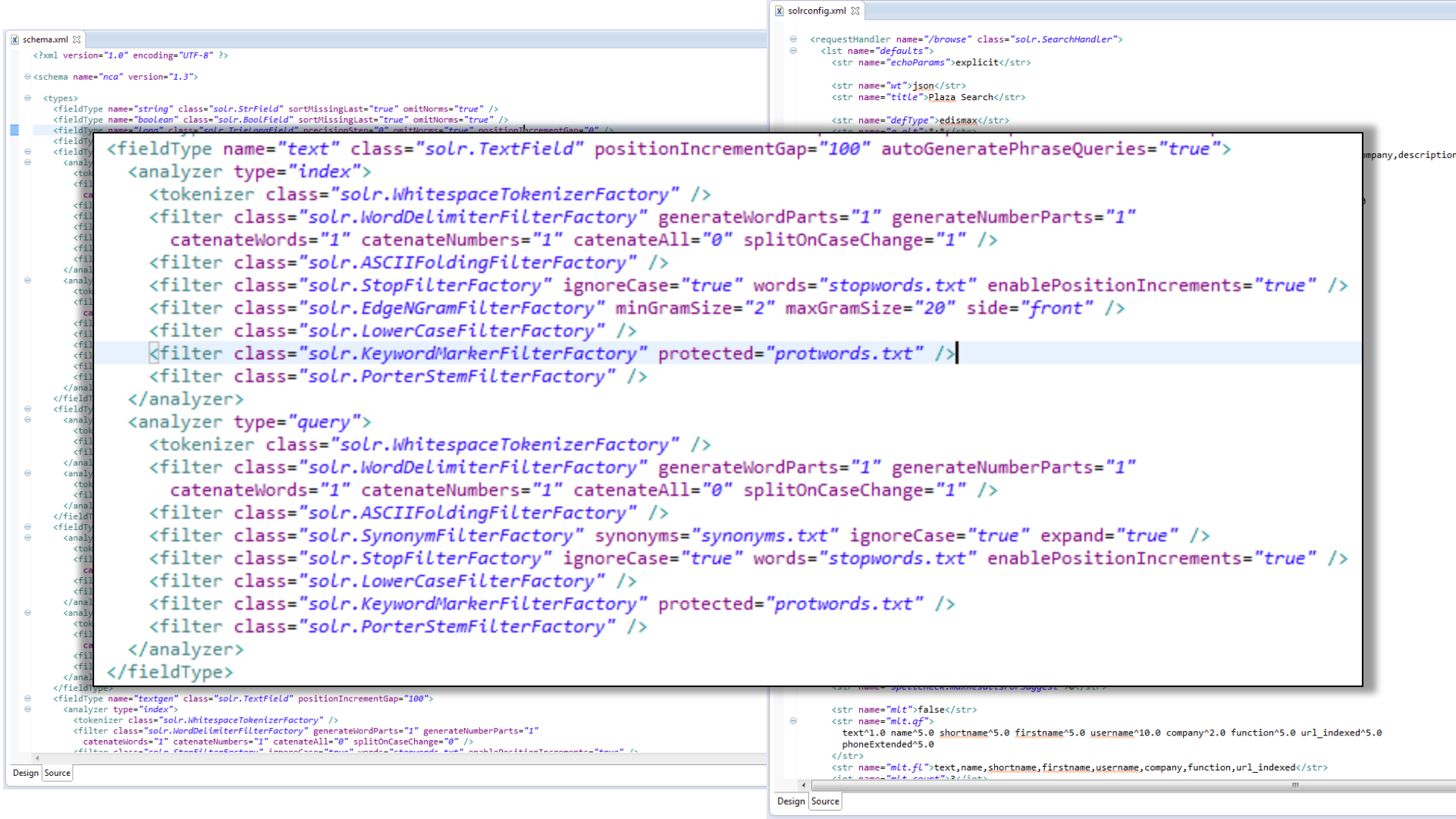
We certainly do not want to build this ourselves

Apache Tika is used for analyzing (file) contents

Also here, we certainly do not want to build this ourselves



Magic



Apache Solr

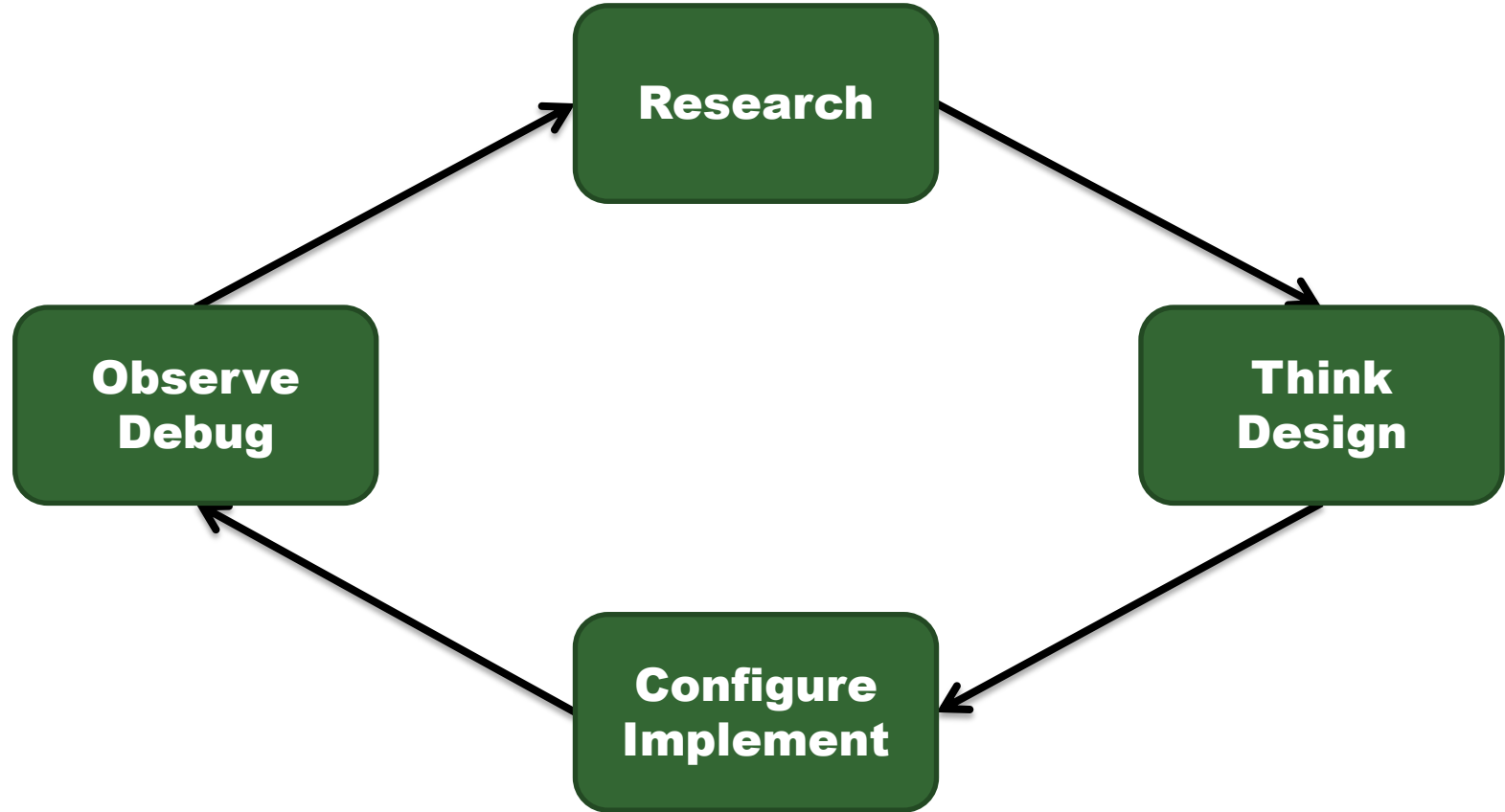
Notes

Apache Solr is a very complex system with a lot of knobs and dials
Most things just seem like magic at the beginning ... or they just do not work

Apache Solr is extremely powerful with a lot of features
You have to know how to configure the features
Most features need a bit more configuration than just a check box for activating it

Configuration options seem very confusing at the beginning
You do not need to understand everything from the start
Defaults are relatively sensible and the example applications are good starting point

Development Process



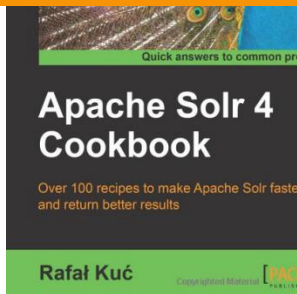
Development Process

Notes

In our experience, Apache Solr works best with a very iterative process
Definition of Done is very difficult to specify for search use cases

Iterate through:

- Researching
- Thinking / Designing
- Implementation / Configuration / Testing
- Observing / Analyzing / Debugging



Solr Wiki [Login](#)

FrontPage

FrontPage

RecentChanges

FindPage

HelpContents

Immutable Page

Info

Attachments

More Actions:



Apache

Solr



StackExchange v

[sign up](#) [log in](#) [careers 2.0](#)



stackoverflow

Questions

Tags

Tour

Users

Tagged Questions

info

newest

6

featured

frequent

votes

active

unanswered

Welcome to

Apache Solr is an open source search server based on the Lucene Java search library.

[learn more...](#) | [top users](#) | [synonyms \(1\)](#)

Contents

1. [Welcome to the Apache Solr Wiki](#)
2. [Solr Documentation](#)
3. [General Information](#)
4. [Solr Development](#)
5. [Using Solr](#)
 1. [Installation and Configuration](#)
 2. [Search and Indexing](#)
 3. [SolrCloud](#)
 4. [Advanced Tools](#)
 5. [Tips, Tricks and Troubleshooting](#)
 6. [Solr Clients](#)
 7. [Operations and Performance](#)
 8. [User-contributed Content](#)
6. [How to edit this Wiki](#)

-1

votes

0

answers

[Generic query builder for nested proximity queries in SOLR](#)

I am trying to build a Java program to transform user entered search criteria to SOLR compatible queries. For example user might enter - ((company or executive) NEAR(5) (retain or hire)) NEAR(5) ...

[solr](#)

4 views

asked 1 hour ago



[g2293](#) [g2293](#)

1

0

votes

0

answers

[solr join function to query documents in multiple cores](#)

I use solr join to query documents from two cores, my cores is defined as follows: Post core: <fields> <!-- general --> <field ...

[solr](#)

[lucene](#)

3 views

asked 5 hours ago



[bright](#)

30 ● 11

Research



Dashboard

Logging

Core Admin

Java Properties

Thread Dump

Core Selector

Instance

Start 10 days ago

Versions

solr-spec 4.5.0
solr-impl 4.5.0 1527178 - jpountz - 2013-09-28 14:11:12
lucene-spec 4.5.0
lucene-impl 4.5.0 1527178 - jpountz - 2013-09-28 14:05:42

JVM

Runtime Oracle Corporation Java HotSpot(TM) 64-Bit Server VM (1...

Processors 12

Args
-Djava.io.tmpdir=/wwwprod/plaza//temp
-Dcatalina.home=/opt/tomcat6
-Dcatalina.base=/wwwprod/plaza/
-Djava.endorsed.dirs=/opt/tomcat6/endorsed
-Djava.util.logging.manager=org.apache.juli.ClassLoader...
-Dcom.sun.management.jmxremote.ssl=false
-Dcom.sun.management.jmxremote.password.file=/www...
-Dcom.sun.management.jmxremote.port=4043
-Dconfluence.browser.language.enabled=false

System

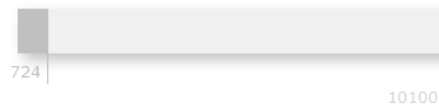
Physical Memory 99.0%



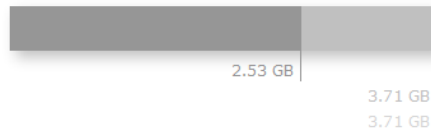
Swap Space 44.4%



File Descriptor Count 7.2%



JVM-Memory 68.1%



Observe
Debug

Solr Admin Interface

Notes

Apache Solr has a pretty good admin interface

Very helpful for analysis and (manual) monitoring

If you are not familiar with the Solr Admin interface, you should be

Other tools like profilers, memory analyzers, monitoring tools etc. are also useful

Our Requirements

Correctness

Results that match query

Relevance

Results that matter

Speed

"Instant" results

Intelligence

Do you know what I mean?

Intelligence

Do you know what I mean?

synonyms.txt

stopwords.txt

protwords.txt

Solr Configuration Files

Notes

Solr has pretty much out-of-the-box support for stop words, protected works and synonyms

These features look very simple, but they are very powerful

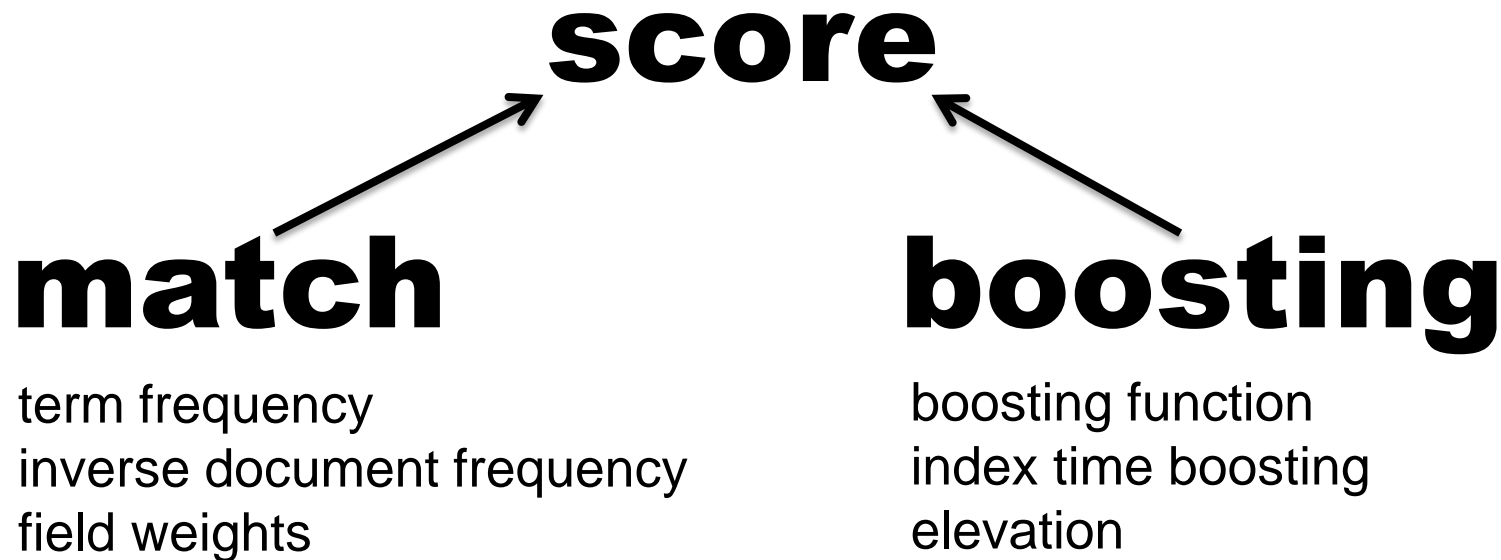
Unless you have a very general search use case, the defaults are **not enough**

Definitely worth **developing a configuration specific to your domain**

Iterate; consider these features for **ranking optimizations** as well

Relevance

Results that matter



Ranking in Solr (simplified)

Notes

Solr determines a score for the results of a query

Score can be used for sorting the results

Score is the product of different factors:

A **query-specific part**, let's call it the **match** value that is computed using
term frequency (tf)

inverse document frequency (idf)

There are also other parameters that can influence it (term weights, field weights, ...)

The match basically says how well a document matches the query

Ranking in Solr (simplified)

A **generic part** (not query specific), let's call this a **boosting** value

Basically represents the general importance that you assign to a document

Includes a ranking function, e.g. based on the age of the document

Includes a boosting value, that is determined at index time (index-time boosting)

We calculate the boost value based on different attributes of the document, such as

- type of resource (people are more important than files)

- status of the project that a document is associated with (closed projects are less important than running projects)

- archive flag (archived resources are less important)

- ...

Ranking Function

`recip(ms(NOW,datestamp),3.16e-11,1,1)`

Index-Time Boosting

Regression **Ranking Testing**

```
assertRank("jira", "url",  
           "https://extranet.netcetera.biz/jira/", 1);  
assertRank("jira", "url",  
           "https://plaza.netcetera.com/.../themas/JIRA", 2);
```

Regression Testing for the Ranking

Notes

Ranking is influenced by various factors

We have continuously executed tests for the ranking

Find ranking regressions as soon as possible

Tests are executed every night, not just with code changes

War Stories

War Story #1:

Disk Space

Situation

Notes

Search is often extremely slow, response times of 20-30s

Situation improves without any intervention

Problem shows up again very soon

Other applications in the same Tomcat server are brought to a grinding halt

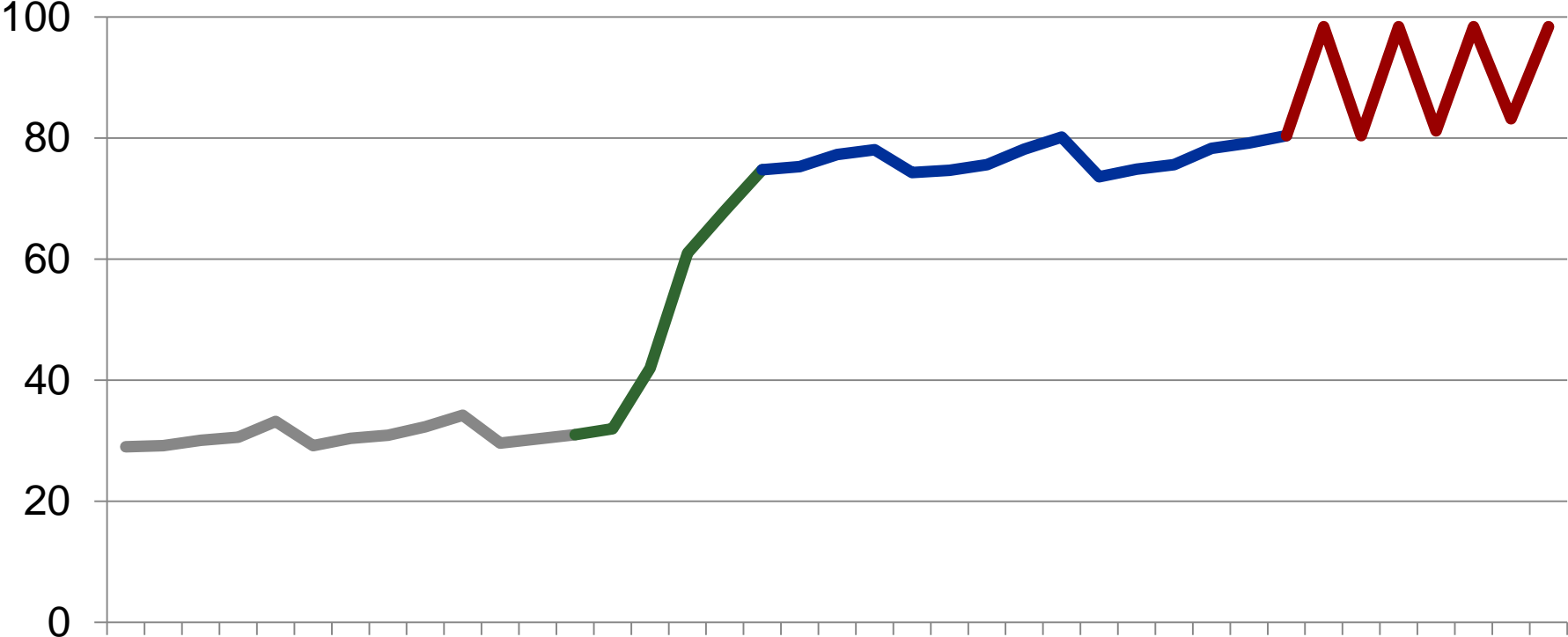
No releases within the last 7 days

No significant data changes in the last 7 days

2-3 weeks earlier a new data sources have been added

Index had grown by a factor of 2, but everything worked fine since then

Disk Usage (fake diagram)



Lucene Index – Disk Usage

Notes

Index needs optimization from time to time when you update it continuously

Index optimization uses a lot of resources, i.e. CPU, memory and disk space

Optimization requires twice the disk space than the optimal index

When your normal index uses > 50% of the available disk space, it's already too late

It's difficult to get out of this situation (without adding disk space)

Deleting stuff from the index does not help, as you need an optimization

Lessons Learned

We need least **2-3 times** as much space as the "ideal" index needs

If your index has grown **beyond 50%**, it's already **too late**.

Disk Usage Monitoring has to be improved

Some problems take a **long time to show themselves**

Testing long-term effects and **continuous delivery** clash to some extent

War Story #2:

Free Memory

Situation

Notes

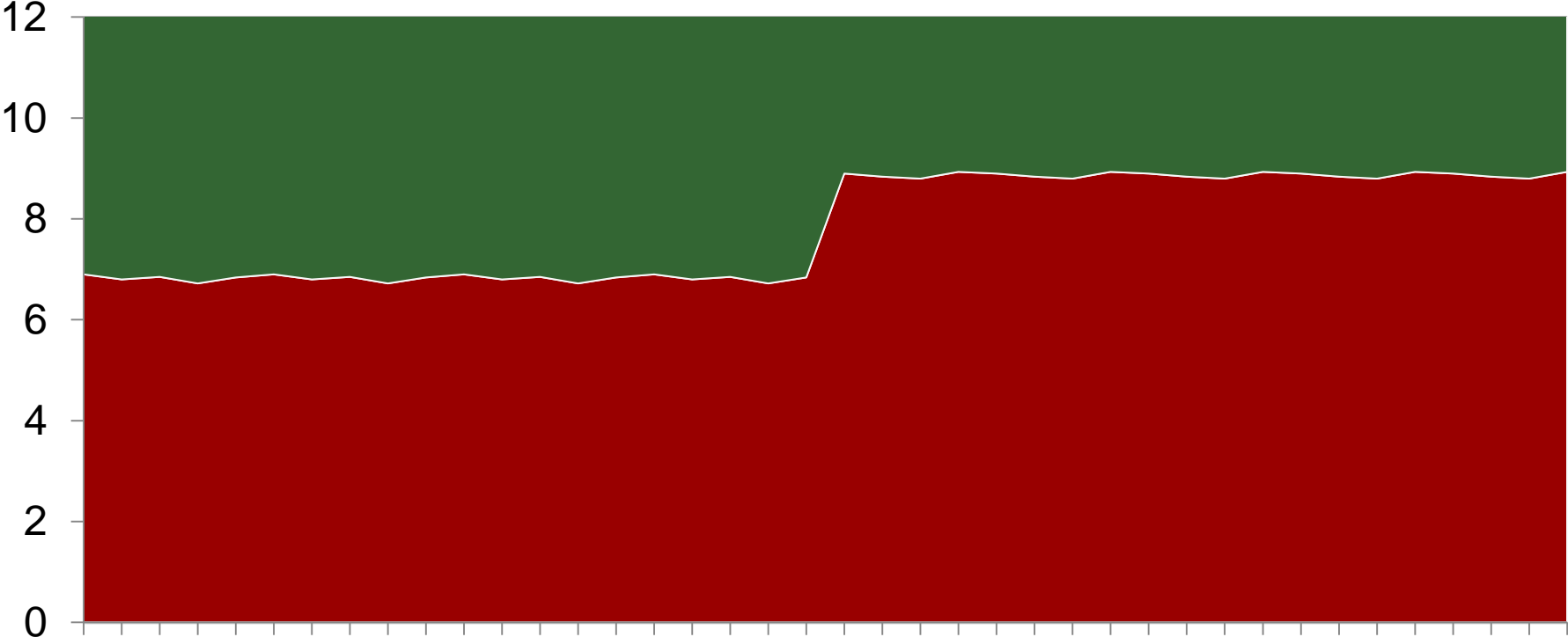
Search is always extremely slow, response times of 20-30s

Other applications in the same Tomcat server show normal performance

No releases within the last few days

No significant data changes in the last few days

Memory Usage (fake diagram)



I/O Caching

Notes

OS uses "free" memory for caching

I/O caching has a HUGE impact on I/O heavy applications

Solr (actually Lucene) is a I/O heavy application

Lessons Learned

Free memory != **unused** memory

Increasing the **heap size** can **slow down Solr**

OS does a better job at caching Solr data than Solr

War Story #3:

Know Your Maths

Situation

Notes

Search starts up very fine and is reasonably fast

Out Of Memory Errors after a couple of hours

Restart brings everything back to normal

Out Of Memory Errors come back after a certain time (no obvious pattern)

Analysis

Notes

Analysis of the memory usage using heap dumps

Solr Caches use up a lot of memory (not surprisingly)

Document cache with up to 2048 entries

Entries are dominated by the content field

Content is limited to 50 KByte by the indexers (or so I thought)

Content abbreviation had a bug

Instead of the 50KByte limit of the indexer, the 2 MByte limit of Solr was used

$2048 * 2 \text{ MByte} = 4\text{GByte}$ for the document cache

Heap size at that time = 4GByte

Lessons Learned

Heap dumps are your friends

Study your heap from time to time, even if you do not have a problem (yet)

Test your **limiters**

War Story #4:

Expensive Features

Situation

Notes

Search has become slower and slower

We added a lot of data, so that's not really surprising

Analysis into different tuning parameters

Analysis into the cost of different features

Highlighting

Plaza Infostore Zimbra Mailstore JIRA Extranet TitraKK Tickets More ▾

Plaza Search

netcetera | Plaza

logging

Search Tips

3743 results for logging

Everything

Wiki

Files

Emails

Web

Issues

People

Companies

Projects

Topics

Resources

Books & Magazines

70% of the response time

3 Introduction NOTE: we might mainly refer to Girders and write here as little as possible. **Logging** in General Rules for **Log** Messages There are a few points to be obeyed when making **log** messages: Performance. **Logging** is usually not a performance issue, except if **logs** statements are generated a few thousand times a second or the generation of the **log** message itself is expensive (e.g. complex toString operations). Production use. **Logging** is not only something to be activated during

Space: **Themas**

logging

<https://plaza.netcetera.com/wiki/display/swc0422/logging>
2013-03-19 11:18 | 3.23 KB | Wiki | HTML | Rolf Koch, Wolfgang Habicht

Logging OBT Coop On swcpub2 **logs** to falsenone/wwwprod/swc/**logs**/obtcoop/obtcoop. **log** will not be rotated. **Logs** trace to catalina.out during

Lessons Learned

Some features are cool, but also **very expensive**

Think about what you need to **index** and what you need to **store**

Consider loading stuff "**offline**" and **asynchronously**

Consider loading stuff from **other data sources**

A few words on
Scaling

Solr Cloud – Horizontal and Vertical Scaling

Support for Replication and Sharding

Added with Apache Solr 4

Based on Apache Zookeeper

Replication

Fault tolerance, failover

Handling huge amounts of traffic

Sharding

Dealing with huge amounts of data



Geographical Replication

Geographical Replication

Notes

Load is not an issue for us, but response time is

We have multiple geographically distribute sites

Network latency is a big factor of the response time if you are at a "far away" location

We have been thinking of setting up replicas of the search engine at the different locations

Relevance-Aware Sharding

Relevance-Aware Sharding

Notes

Normal sharding distributes data on different, but equal nodes

We have been thinking about creating deliberately different nodes for the distribution of the data:

Node 1

- extremely fast
- small index, i.e. small amount of data
- lots of memory, CPU, really fast disks

Node 2

- lots more data
- big, but slower disks
- less memory and CPU

Frontend would send queries to both nodes and show results as they come in.

Distribution of the data would be based on the (query independent) boosting value

Wrapping Up

Search rocks

Apache Solr rocks

Learning Curve

Definition of **Done**

Continuous **Inspection**
Continuous **Improvement**

Get your hands dirty

Ranking Optimizations

Continuous **Testing** and **Monitoring** for **Ranking** and **Performance Issues**

Verification
of features can take
a **long time**

Cool **side projects** rock

Contact



Corsin Decurtins

corsin.decurtins@netcetera.com

[@corsin](#)

References

Notes

Apache Solr

<http://lucene.apache.org/solr/>

Apache Solr Wiki

<http://wiki.apache.org/solr/>

Apache Solr on Stack Overflow

<http://stackoverflow.com/questions/tagged/solr>