



Lucene: Generic Data Indexing

Mike Cannon-Brookes

CEO, Atlassian Software Systems

Java Champion





WARNING

**The following presentation
contains egregious product
placement and lots of text.**



Indexing:

"The process of converting a collection of data into a format suitable for easy search and retrieval."



JIRA: Issue Search

JIRA User: Mike Cannon-Brookes [Filters](#) | [Profile](#) | [Log Out](#)

[HOME](#) [BROWSE PROJECT](#) [FIND ISSUES](#) [CREATE NEW ISSUE](#) [ADMINISTRATION](#) QUICK SEARCH:

Filter: [View](#) [Edit](#) [New](#) [Manage](#)

You are currently using a new, unsaved search.
 [Save](#) it as a filter

[View & Hide](#) [View >>](#)

Project: Bamboo
BNP Consulting
Citigroup Change Request
Cohotion
Confluence

Issue Type: Any
Standard Issue Types
 Bug
 Improvement
 New Feature
 Support Request

Components / Versions

Fix For: ~~2.0 Fix version~~
Unreleased Versions
2.2.11
2.3.4
2.3.x
2.4.3

Components: Any
No Component
Administration
Application Server Support
Attachments

Affects Versions: Any
No Version
Released Versions
2.4.2
2.4.1
2.4

Text Search

Query:

Query Fields: Summary Description
 Comments Environment

Issue Navigator

Displaying issues 1 to 10 of 10 matching issues. [\[Permlink \]](#)

Current View:
Browser ([Current Fields](#) | [Printable](#) | [Full Content](#)) | [XML](#) | [RSS](#) ([Issues](#) | [Comments](#)) | [Word](#) | [Excel](#) ([All fields](#) | [Current fields](#)) | [Charts](#)

Bulk Change: [all 10 issue\(s\)](#)
 [Configure](#) your Issue Navigator

T	Key	Summary	Reporter	Assignee	Status	Res	Updated
	CONF-8078	Error page if new password doesn't match Crowd password validation	Matt Ryall	Unassigned	Open	UNRESOLVED	19/Mar/07
	CONF-8055	Document process for moving from evaluation to commercial cluster license	Matt Ryall	Unassigned	Open	UNRESOLVED	14/Mar/07
	CONF-8050	zip_src from tiny mce served without caching headers on extranet	Scott Farquhar	Matthew Jensen	Reopened	UNRESOLVED	19/Mar/07
	CONF-7987	\$(baseUrl) is not being substituted correctly in daily update email.	Rob Di Marco	Unassigned	Needs Verification	UNRESOLVED	12/Mar/07
	CONF-7975	Migrate c.a.c to hibernate user repository	Christopher Owen	Unassigned	Open	UNRESOLVED	12/Mar/07
	CONF-7893	Link to Plugin Repository is 'plugin.repository.link'	David Soul	Unassigned	Open	UNRESOLVED	12/Mar/07
	CONF-7774	Global Activity link on Space Activity page is not correct	Agnes Ro	Christopher Owen	Open	UNRESOLVED	18/Mar/07
	CONF-7456	Intermittent problem loading stylesheets and javascript on extranet	Scott Farquhar	Christopher Owen	In Progress	UNRESOLVED	19/Mar/07
	CONF-6953	RSS parsing and content-type problems	Matt Ryall	Unassigned	Open	UNRESOLVED	12/Mar/07
	CONF-6058	Group picker for page restrictions silently hides non-member groups, may need clarifying sentence added	David Soul	Unassigned	Open	UNRESOLVED	19/Mar/07



JIRA: Lucene History

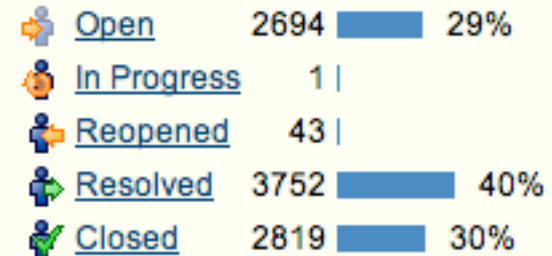
- **1.4 - Use DB for all queries, Lucene only if full text search - results 'merged'**
- **2.0 - Use Lucene for all search, Java for permissioning - results iterated over and non-view stripped**
- **2.2 - Use Lucene for all queries including perms - sorting still done in Java**
- **2.4 - Use Lucene for all queries, retrieving issues and displaying - no DB access at all!**
- **3.0 - Switch "stats" over to using Lucene via HitCollectors**



JIRA: Statistics

All from Lucene!

Project Summary



JIRA Atlassian JIRA
User: Clark Kent | History | Filters

HOME BROWSE PROJECT FIND ISSUES CREATE NEW ISSUE
QUICK SEARCH:

Atlassian JIRA

[Configure:](#)

Statistics: JIRA (Fix For Versions (non-archived))

	3.6.3	20	1%
	3.6.x	136	5%
	3.7	183	7%
	3.8	4	
	Unscheduled	2421	88%

Created vs Resolved Issues: JIRA

Issues: **140** created and **93** resolved
Period: last 30 days (grouped daily)
[View detailed data table >>](#)

Pie Chart: JIRA [more detail >>](#)

Issues: **9382**.
[View detailed data table >>](#)

Statistics: Confluence (Fix For Versions (non-archived))

	2.0	1	
	2.1.5	2	
	2.2	1	
	2.2.2	2	
	2.2.7	15	1%
	2.2.x	288	14%
	2.3	10	
	2.4	4	
	Buffy	4	

Saved Filters (Create New | Manage Filters)

All issues reported by me	1
All Resolved Confluence bugs	1881
All Resolved .JIRA Bugs	2596



Lucene: Full Text Search

- **Text Analysis & Stemming**
 - “Michael jogs in the park” > “michael, jog, park”
- **Proximity Queries**
 - “cat NEAR dog”
- **Wildcard Queries**
 - “jog*”, “j?g”
- **Results returned scored by *relevance***



Lucene: Generic Data Indexing (GDI)

- **Fast retrieval of complex data objects**
 - Built from one database, multiple databases, files, anywhere
 - Not a single table - just use a database index

Issue Attributes

Reporter: Any User

Assignee: Any User

Status: Any
Open
In Progress

Resolutions: Any
Unresolved
Fixed

Priorities: Any
Blocker
Critical



Lucene: Generic Data Indexing (GDI)

- **Powerful pre-built query tools**
 - RangeQuery, BooleanQuery etc

“select issues created between 2001 and 2004, with no components, no versions, still unresolved that have > 4 votes”


The screenshot displays a search filter panel with two main sections:











- ▼ Dates and Times**: This section contains several date-related filters, each with a text input field and a calendar icon to its right:
 - Created After: [input field]
 - Created Before: [input field]
 - Created: From [input field] To [input field] (with a note: "Use this picker for relative dates")
 - Updated After: [input field]
 - Updated Before: [input field]
 - Updated: From [input field] To [input field] (with a note: "Use this picker for relative dates")
 - Due After: [input field]
 - Due Before: [input field]
 - Due Date: From [input field] To [input field] (with a note: "Use this picker for relative dates")
- ▼ Actual vs Estimated Work Ratio**: This section contains a filter for work ratio limits:
 - % Limits: Min [input field] Max [input field] (with a help icon) (with a note: "Enter a minimum, maximum or range limit")



Lucene: Generic Data Indexing (GDI)

- Results returned sorted in *custom order*
 - Sort, SortField

Key 	Summary	
CONF-8078	Error page if new password doesn't match Crowd password validation	F
CONF-8055	Document process for moving from evaluation to commercial cluster license	M
CONF-7987	`\${baseUrl}` is not being substituted correctly in daily update email.	F
CONF-7975	Migrate c.a.c to hibernate user repository	C

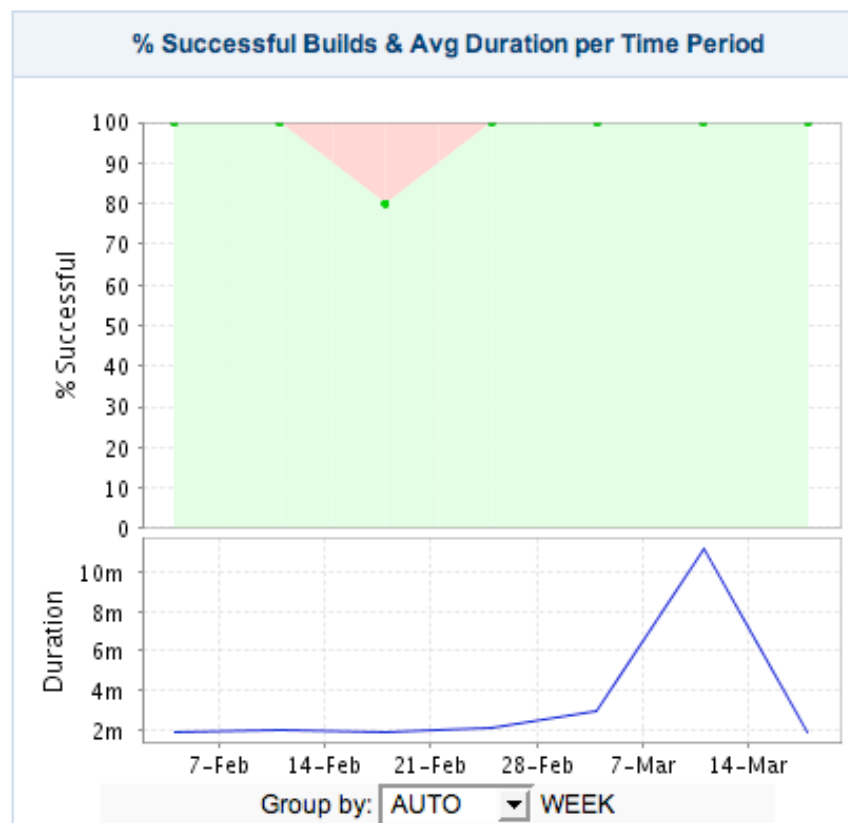
Versions	
<input type="checkbox"/> Manage versions	(displayed in the order of newest first)
 Buffy	
 2.5	
 2.4.x	
 2.4.3	
 2.4.2	12/Mar/07
 2.4.1	12/Mar/07
 2.4	25/Feb/07
 2.3.x	
 2.3.4	
 2.3.3	14/Feb/07



Lucene: Generic Data Indexing (GDI)

- AOP-like result filtering and hit collection
 - QueryFilter and HitCollector

		Pages					New
	View	Create	Export	Restrict	Remove	Create	R
Eugene Katz	✓	✓	✓	✓	✓	✓	
Ivo Verlaek	✓	✓	✗	✗	✗	✗	
William Anderson	✓	✓	✓	✓	✓	✓	✓





Lucene: Generic Data Indexing (GDI)

- Integrated full text search - only *if* you need it!
- “Free!”

select issues created between 2001 and 2004,
with no components, no versions,
still unresolved that have > 4 votes
and match the query “dash”*



Database V1

Issues

Issue	Summary	Assignee	Reporter
JRA-1	Buy milk	Fred	
JRA-2	Collect laundry	Bill	Fred

Query: `select * from issues where assignee = 'fred'`



Database V2

Issues

Issue	Summary
JRA-1	Buy milk
JRA-2	Collect laundry

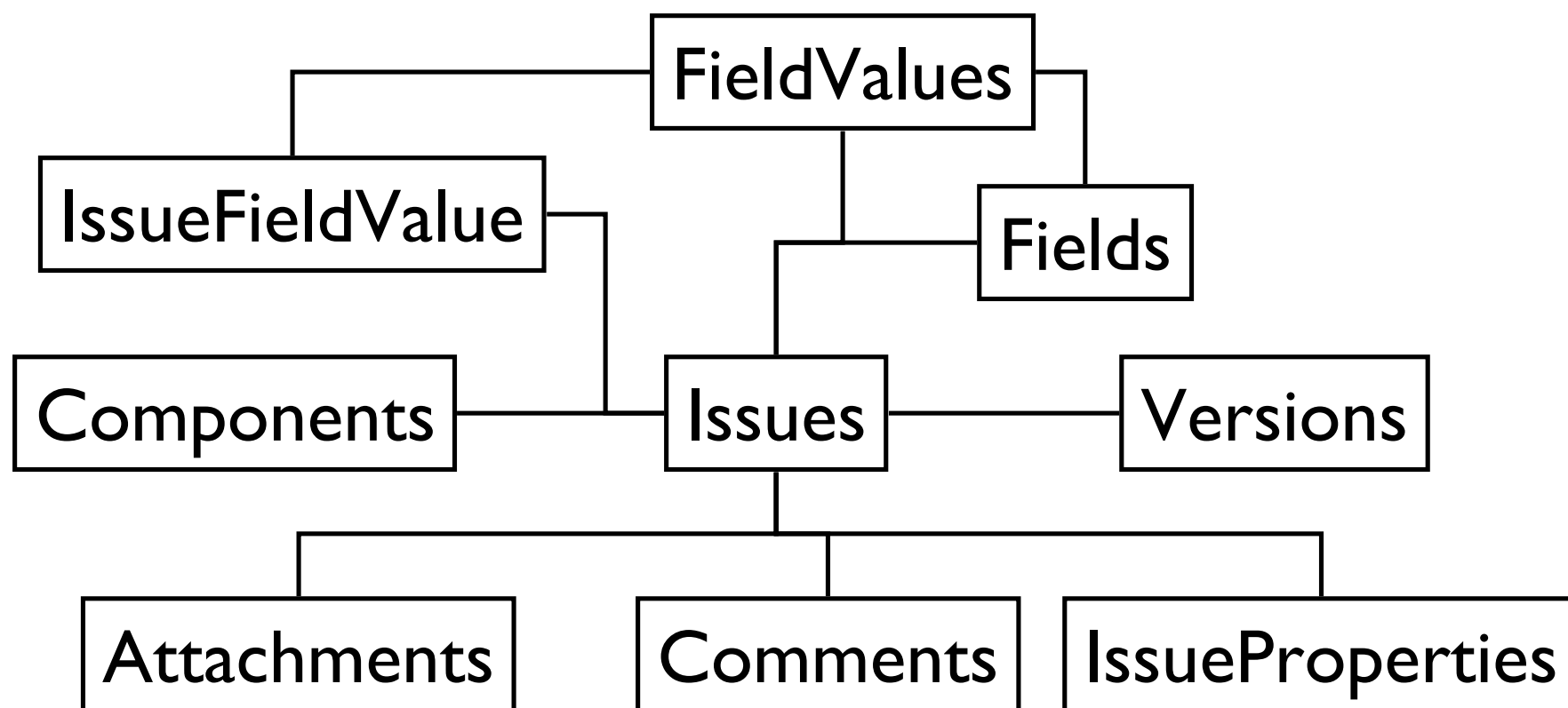
Fields

Issue	Field	Value
JRA-1	Assignee	Fred
JRA-2	Assignee	Bill
JRA-2	Reporter	Fred

Query: select * from issues, fields
where fields.field = 'Assignee' and fields.value = 'fred'
And fields.issue = issues.issue



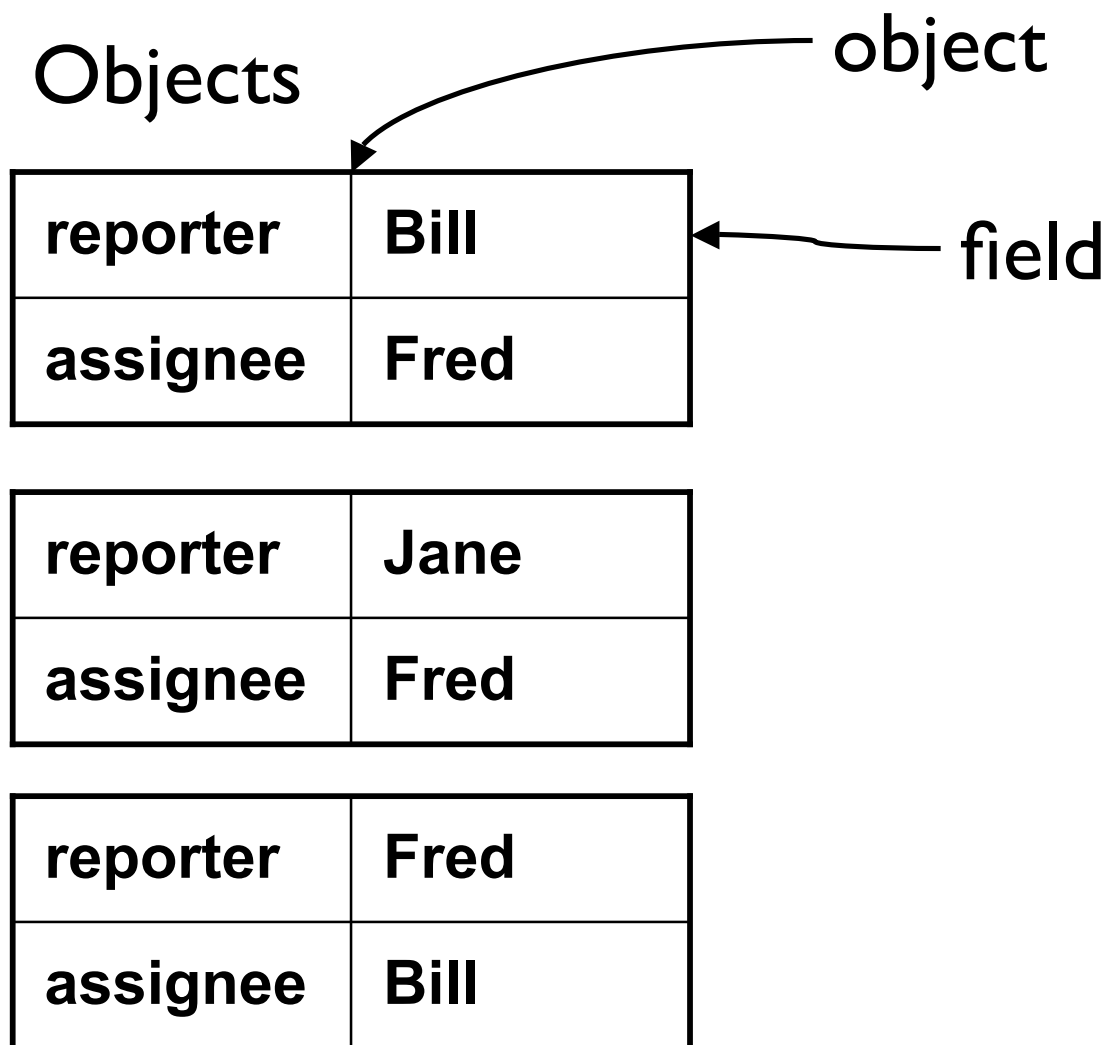
Database V3



Query: ???



How is Lucene fast?





How is Lucene fast?

Objects

reporter	Bill
assignee	Fred

reporter	Jane
assignee	Fred

reporter	Fred
assignee	Bill

Docs

docId: 1



docId: 2



docId: 3



document

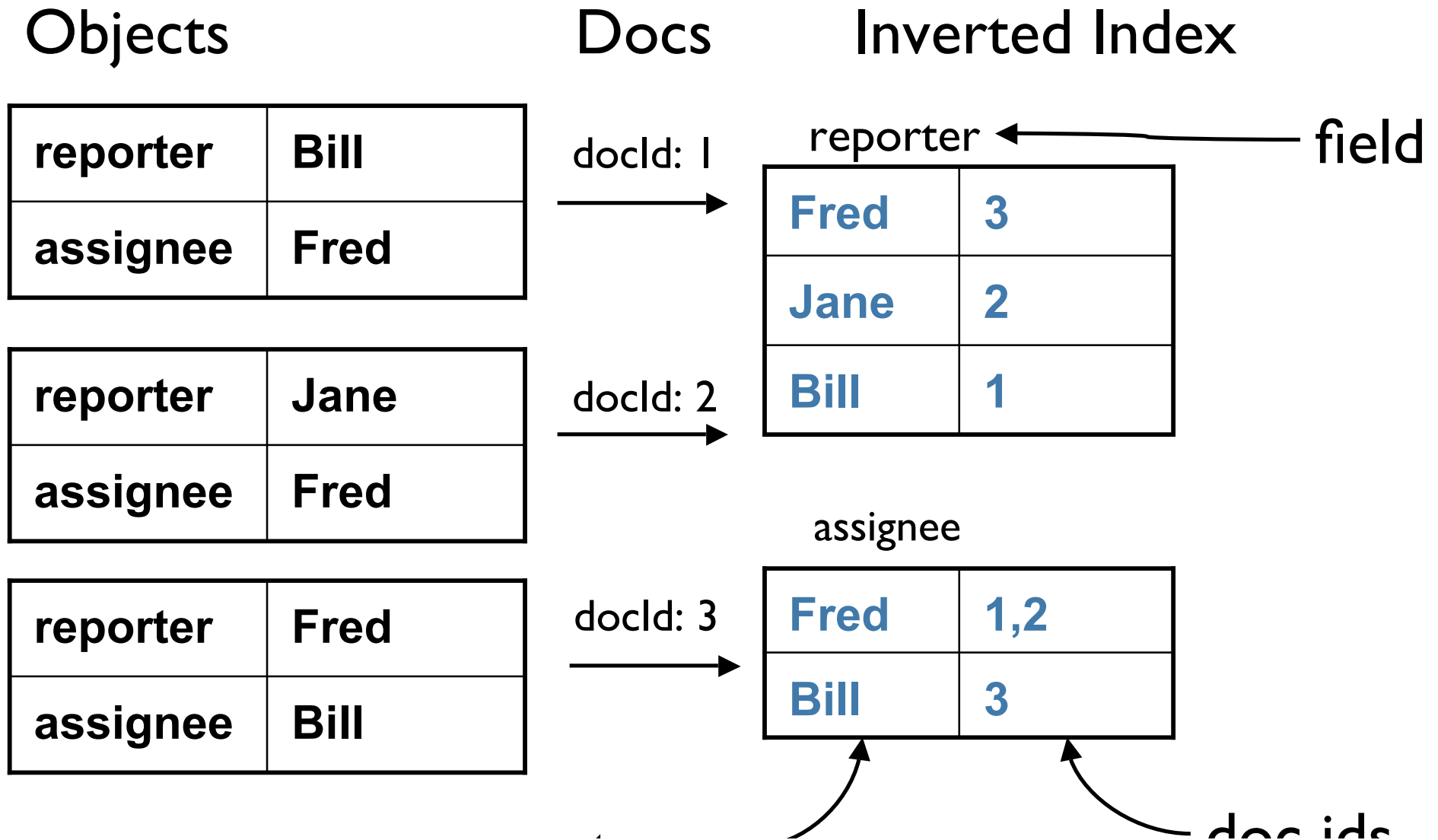


```
reporter : Bill  
assignee : Fred  
component : 1  
component : 4  
component : 5  
created : 20070320
```

```
field : value ...
```



How is Lucene fast?





How is Lucene fast?

Objects

reporter	Bill
assignee	Fred

reporter	Jane
assignee	Fred

reporter	Fred
assignee	Bill

Docs

docId: 1
→

docId: 2
→

docId: 3
→

Inverted Index

bitset

reporter

Fred	3	001
Jane	2	010
Bill	1	100

assignee

Fred	1,2	110
Bill	3	001



Advantages For Generic Data Indexing

- **Store Denormalised Data**
 - Issue object, fields - one single Lucene document
 - No record de-duplication needed as per SQL query
- **Native Java API**
 - Useful for things like sorting where a DB can't do it
 - Java specific sort algorithm for issue keys
 - Version sequencing - v. complex to do in DB



Advantages For Generic Data Indexing






- **Constant time & capabilities**
 - Our apps are cross platform, OS, JDK, database
 - Lucene works pretty much the same across all of them unlike SQL
 - Local file system access (most commonly) which is faster than DB as no network time
- **Constant index format**
 - Readable from Java, C, Perl, Ruby etc
- **QueryFilters & HitCollectors!**



HitCollector

- **Call back object for hit collection**
- **Great for statistical operations where content / score is irrelevant**
 - JIRA - StatusHitCollector for 'bucketing'
- **Fast because:**
 - Retrieve only fields you need
 - Minimum number of loops

Project Summary

 Open	2694	<div style="width: 29%;"></div>	29%
 In Progress	1		
 Reopened	43		
 Resolved	3752	<div style="width: 40%;"></div>	40%
 Closed	2819	<div style="width: 30%;"></div>	30%



HitCollector

```
public abstract void collect(int doc,  
                             float score)
```

Called once for every non-zero scoring document, with the document number and its score.

If, for example, an application wished to collect all of the hits for a query in a BitSet, then it might:

```
Searcher searcher = new IndexSearcher(indexReader);  
final BitSet bits = new BitSet(indexReader.maxDoc());  
searcher.search(query, new HitCollector() {  
    public void collect(int doc, float score) {  
        bits.set(doc);  
    }  
});
```




QueryFilter

```
public BitSet bits(IndexReader reader)  
    throws IOException
```

Description copied from class: [Filter](#)

Returns a BitSet with true for documents which should be permitted in search results, and false for those that should not.

Specified by:

[bits](#) in class [Filter](#)

Throws:

[IOException](#)



Atlassian: Examples Of Lucene Usage

- **JIRA - User driven queries an arbitrary data model**
 - Plugins index/search their own 'fields' - future proof!
 - QueryFilters for permissions - cached per request
 - HitCollectors for all statistics / dashboard

The screenshot shows the JIRA Issue Navigator interface. The top navigation bar includes links for HOME, BROWSE PROJECT, FIND ISSUES, CREATE NEW ISSUE, and ADMINISTRATION. The user is identified as Mike Cannon-Brookes. The interface displays a list of issues with columns for Key, Summary, Reporter, Assignee, Status, Res, and Updated. The issues listed are:

T	Key	Summary	Reporter	Assignee	Status	Res	Updated
	CONF-8078	Error page if new password doesn't match Crowd password validation	Matt Ryall	Unassigned	Open	UNRESOLVED	19/Mar/07
	CONF-8055	Document process for moving from evaluation to commercial cluster license	Matt Ryall	Unassigned	Open	UNRESOLVED	14/Mar/07
	CONF-8050	zip_src from tiny mce served without caching headers on extranet	Scott Farquhar	Matthew Jensen	Reopened	UNRESOLVED	19/Mar/07
	CONF-7987	`\${baseUrl}` is not being substituted correctly in daily update email	Rob Di Marco	Unassigned	Needs Verification	UNRESOLVED	12/Mar/07

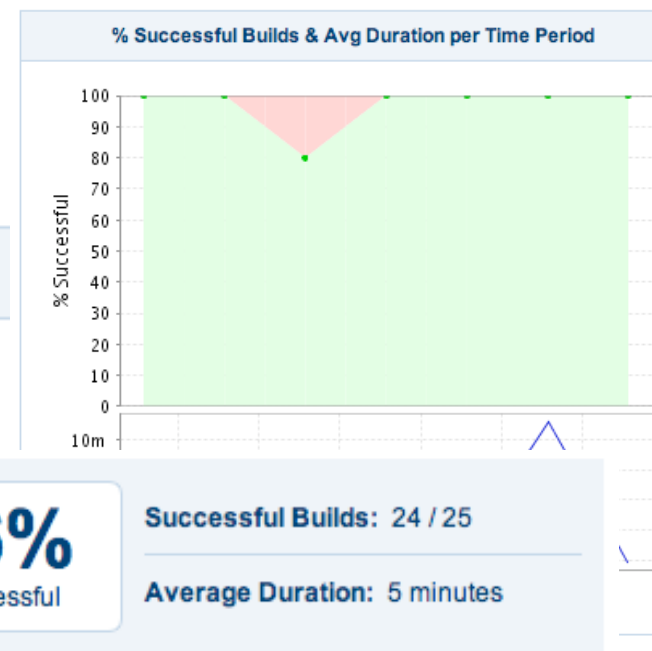


Atlassian: Examples Of Lucene Usage

- **Bamboo - Build telemetry statistics via Lucene**
 - Fast over millions of rows - data on every test/suite/build run, ever.
 - Use Lucene to aggregate data into useful statistics
 - HitCollectors used extensively for telemetry data

Recent Failures

- Average time to fix a failure: **2 days, 19 hours, 16 minutes**
- Average number of builds between fixes: **1 builds**
- The longest time taken to fix a failure is **2 days, 19 hours, 16 minutes**
- The greatest number of builds taken to fix a failure is **1**, from failure sta





Problems For GDI

- **One Big Singleton**
 - Updates require serialization - indexes are write once, read many
 - Jira vs Confluence different access/write strategies
- **Delete / Update Operations**
 - Lucene wasn't built for fast changing data
 - Delete operation is just a flag op & Update requires delete / re-add
 - Fixed in Lucene 2.1 - <http://issues.apache.org/jira/browse/LUCENE-565>
- **Writing Is Expensive**
 - Opening/closing reader/writers proportional to index size



Problems For GDI

- **Timing Of `index.optimise()`**
 - **Indexes get fragmented - `optimise()` defrags**
 - **Tricky to time this as v. slow on large indexes**
 - **Eden space strategy can solve this**
 - **Small index for 'updated' data, large for 'old' data**
 - **Optimize large index rarely, small frequently - like GC.**
 - **MultIndexSearcher allows search on multiple indexes like one**



Problems For GDI

- **Non Transactional**

- DB can have data that index misses, or vice versa
- Compass is a solution here - haven't tested
- Otherwise, architect correct design knowing Lucene

- **Optionitis**

- Write settings can require a lot of knowledge and tuning
- eg MAX_MERGE_DOCS

- **Local storage - can be a problem in a cluster**

- See my other presentation for clustering strategies!



GDI Lucene Usage Models

● JIRA

- Synchronous indexing > tricky locking problems at scale
- More updates than creates > heavier index load
 - Fixed with Lucene 2.1!
- Slower updates, statistics always correct

● Confluence

- Asynchronous indexing
- Updates are queued, flushed every minute
- Clusterable and faster 'net' time for user
- 'Recent Updates', 'Search' up to 1 min inaccurate



Tips

- **Use derived data only, so can be recreated at will**
- **Store -1 for null because nature of fields**
 - **Can't query Lucene for 'lack of a field' - ie "No Component"**
- **Keep open a single searcher and 'flip it' after writing**
- **ThreadLocals are valuable in web apps**
 - **Use for Searchers, BooleanQueries and QueryFilters that are expensive to create per search but don't change per request (10s of queries per request)**
- **Understand Lucene to adjust your usage to your app**
- **Index dates to highest granularity possible, prevent term explosion**
 - **Remember Lucene stores? YYYYMMDD vs YYYYMMDDHHmmSSSS**



Links

- **Luke - useful tool to examine indexes**
 - <http://www.getopt.org/luke>
- **Lucene In Action - *awesome* book**
 - <http://www.lucenebook.com>
- **Compass - Lucene abstraction framework**
 - <http://www.opensymphony.com/compass>



Q & A

P.S. Java guru? Atlassian needs engineers!

- Sydney or San Francisco.

<http://www.atlassian.com/about/jobs>

Email me: mike@atlassian.com

