



# Becoming a content-driven, modular application A Case Study

Brett Porter, MaestroDev

[brett@maestrodev.com](mailto:brett@maestrodev.com)

Produced by



# Brett Porter

# Brett Porter

- Apache Software Foundation
  - Apache Maven, Archiva, Continuum, NPanday, Infrastructure, others
  - Member and Director

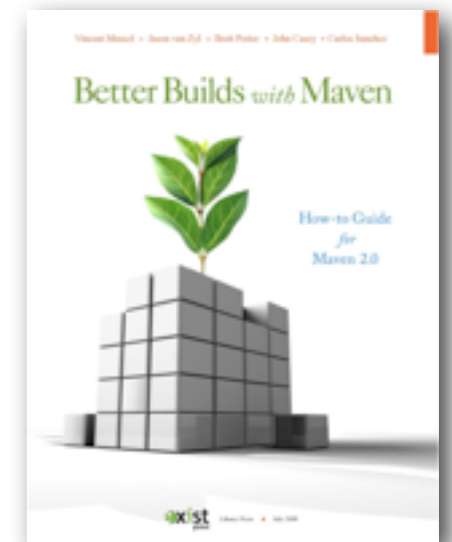
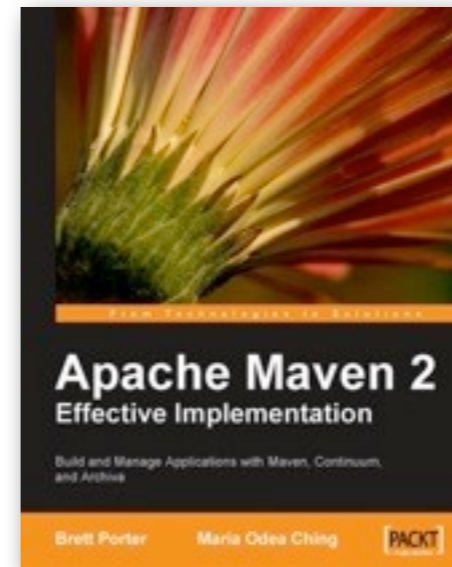
# Brett Porter

- Apache Software Foundation
  - Apache Maven, Archiva, Continuum, NPanday, Infrastructure, others
  - Member and Director
- MaestroDev
  - CTO
  - Directing Maestro development



# Brett Porter

- Apache Software Foundation
  - Apache Maven, Archiva, Continuum, NPanday, Infrastructure, others
  - Member and Director
- MaestroDev
  - CTO
  - Directing Maestro development
- Co-author
  - Apache Maven 2: Effective Implementation
  - Better Builds with Maven



# Content-related Experience

# Content-related Experience

- This page left intentionally blank

# A Practical Example

- This is about a change we made in our project
- May apply to you if you have a similar challenge
- Only certain types of applications fit as a content application



# A Practical Example

- This is about a change we made in our project
- May apply to you if you have a similar challenge
- Only certain types of applications fit as a content application



- Apache Archiva
  - You can check out the code for yourself, it's open source
  - <http://svn.apache.org/repos/asf/archiva/trunk/>



# Archiva: Some Background

- Repository Manager for Maven (and other similar tools)
- Naturally hierarchical content based on the Maven repository format
- Basically an artifact file server with a custom interface
  - rule-based retrieval and management of artifacts and associated metadata
  - access directly over HTTP and WebDAV, a user-driven web interface, and some web services
  - artifacts are typically binaries, ranging from small to multi-gigabyte, with information attached from the Maven POM or other sources
  - typically a large number of files, and rapid turnover as new are added and older development snapshots are purged

# Archiva: Some Background

- Repository Manager for Maven (and other similar tools)

- Naturally hierarchical

- Basically an artifact

- rule-based retrieval

- access directly via  
web services

- artifacts are typed  
information attached

- typically a large  
development space

The screenshot displays the Apache Archiva web interface. On the left is a navigation sidebar with sections: Find (Search, Find Artifact, Browse), Manage (Reports, User Management, User Roles, Appearance, Upload Artifact, Delete Artifact), and Administration (Repository Groups, Repositories, Proxy Connectors, Legacy Support, Network Proxies, Repository Scanning, Database). The main content area shows the details for the artifact 'Lang' (commons-lang / commons-lang / 2.3). It includes a description: 'Commons.Lang, a package of Java utility classes for the classes that are in java.lang's hierarchy, or are considered to be so standard as to justify existence in java.lang.' Below this are metadata fields: Repository (releases), Group ID (commons-lang), Artifact ID (commons-lang), Version (2.3), and Packaging (jar). A 'Downloads' box shows 245,274 Jars and 11,115 Poms. A 'POM Snippet' shows the XML dependency declaration. 'Other Details' includes URL, Organisation (The Apache Software Foundation), License (The Apache Software License, Version 2.0), Issue Tracker, and Continuous Integration. 'SCM' information includes Connection and Viewer URLs. The footer shows 'Apache Archiva 1.2-SNAPSHOT' and 'Copyright © 2005-2009 The Apache Software Foundation'.

some

older

# Archiva: Some Background

- Repository Manager for Maven (and other similar tools)

- Naturally hierarchical

- Basically an artifact

- rule-based retrieval

- access to

web services

- artifacts are typed  
information attached

- typically a large  
development space

The screenshot shows the Apache Archiva web interface. At the top, the Archiva logo is on the left, and the current user 'admin (admin)' is on the right. The main content area displays the details for the artifact 'commons-lang / commons-lang / 2.3'. It includes a search bar, a 'Downloads' box showing 245,274 Jar downloads and 11,115 Pom downloads, and a 'POM Snippet' section with the following XML code:

```
<dependency>
  <groupId>commons-lang</groupId>
  <artifactId>commons-lang</artifactId>
  <version>2.3</version>
</dependency>
```

Below the POM snippet, there are sections for 'Other Details' (URL, Organisation, License) and 'SCM' (Connection, Viewer). The footer of the page reads 'Apache Archiva 1.2-SNAPSHOT' and 'Copyright © 2005-2009 The Apache Software Foundation'.

<http://vmbuild.apache.org/archiva/repository/central-proxy/commons-lang/commons-lang/2.3/commons-lang-2.3.jar>

older

# Other Repository Managers

# Other Repository Managers



- Uses Lucene and flat files for metadata
- Has an established anti-database stance
- Focuses on “self-healing” metadata to ensure integrity

# Other Repository Managers



- Uses Lucene and flat files for metadata
- Has an established anti-database stance
- Focuses on “self-healing” metadata to ensure integrity
- Initially used JCR to store everything, including binary artifact data
- Claimed benefits of integrity
- Had reputation for wedging the database
- Harder to import/export the content
- Now seems to support a filesystem-only repository



# Archiva History

- Around in part since March 2005
  - converting Maven 1 to Maven 2 repositories
  - relied heavily on scanning the repository on the filesystem and pulling out metadata
- Grew into a repository manager application as a Maven subproject
  - promoted to top level project Apache Archiva in March 2008
- Architecture was using Lucene as a “database”, but stored everything in its original form

# Archiva 1.0 Architecture

- Leap forward in functionality
- All of storage was re-done, partly using database
  - to be able to query easily and use persistence APIs
  - to do two phase scanning

# This Didn't Work Out So Well

- We used JDO 2 (JPOX) - not a great deal of resources for it
- Two-phase scanning could get out of sync
- Fell into the classic trap - only one person knew how it worked
- A few problems started to crop up
  - database exceptions deep in the stack that were hard to deal with
  - performance concerns
  - memory consumption (particularly with embedded database)
  - lack of extensibility for metadata
  - configuration for initial set up was not necessarily out of the box

# This Didn't Work Out So Well

The image shows a screenshot of the Archiva MRM (Maven Repository Manager) interface. On the left, a list of issues is displayed, each with a title and a status icon (a red arrow pointing up and a green checkmark). The issues are:

- MRM-914** ...in column "DESCRIPTION" that has maximum length of 8192. Please correct your data!
- MRM-951** "Unable to find project model" although everything seems to be perfectly fine
- MRM-729** [MySQL] Specified key was too long; max key length is 765 bytes - in reback
- MRM-657** 'ORA-00910: specified length too long for its datatype' Error when clicking on searched artifact.
- MRM-568** Error in 'update-db-project' consumer during database scanning when a repo that has been removed was re-added again
- MRM-735** Database on MS SQL 2000/2005 fail to be created due to too column length
- MRM-721** ConsumerException when scanning database
- MRM-990** Archiva hangs with connection pool error
- MRM-705** database scanning should be able to be run after repository scanning
- MRM-1001** Allow for easier database upgrade
- MRM-1235** Upgrade task for altering the database schema for the changes in the length of URL fields

On the right, a table shows the dates for a specific issue:

Dates	
Created:	24/Oct/07 6:04 AM
Updated:	23/Feb/10 2:25 AM
Resolved:	23/Feb/10 2:25 AM

# Motivation for Change

- The architecture was holding it back
- Wanted to implement extensible metadata for artifacts

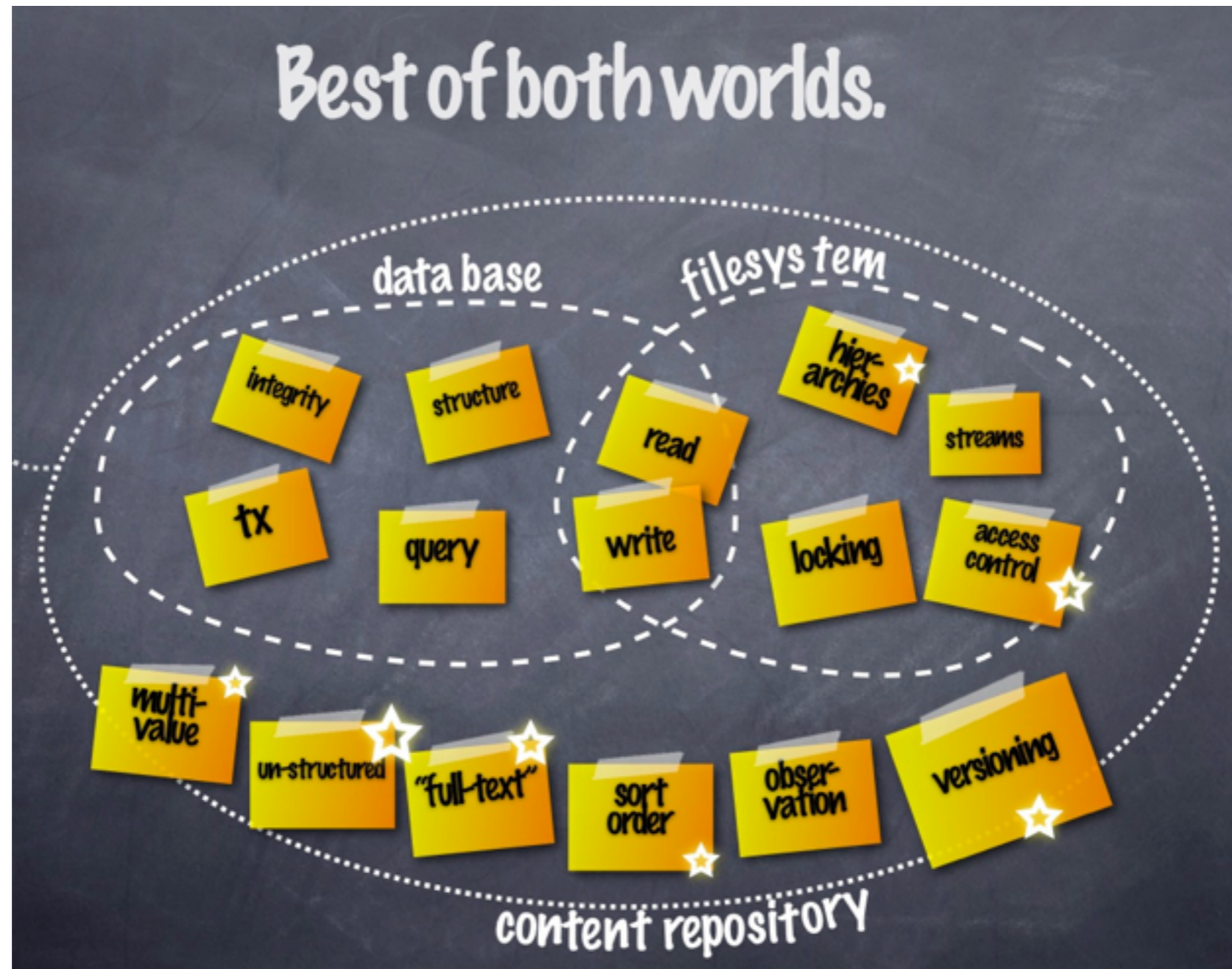
# Back to the Future

- Blend up the strengths of the original architecture, the newer feature set, and the added metadata
- Allow use of unstructured data
- Improve or solve the problems we'd seen
- Remove the database altogether
- Separate the metadata from the storage
- Pick up other improvements on the way, like lazy-loading content inside artifacts and proxying remote repositories
- Move toward a defined target architecture, and keep it working along the way
- Make the application modular: reusable and extensible
- Add a plugin architecture



# Content vs. Database

- <http://java.dzone.com/articles/java-content-repository-best>





# Hierarchical vs. Relational

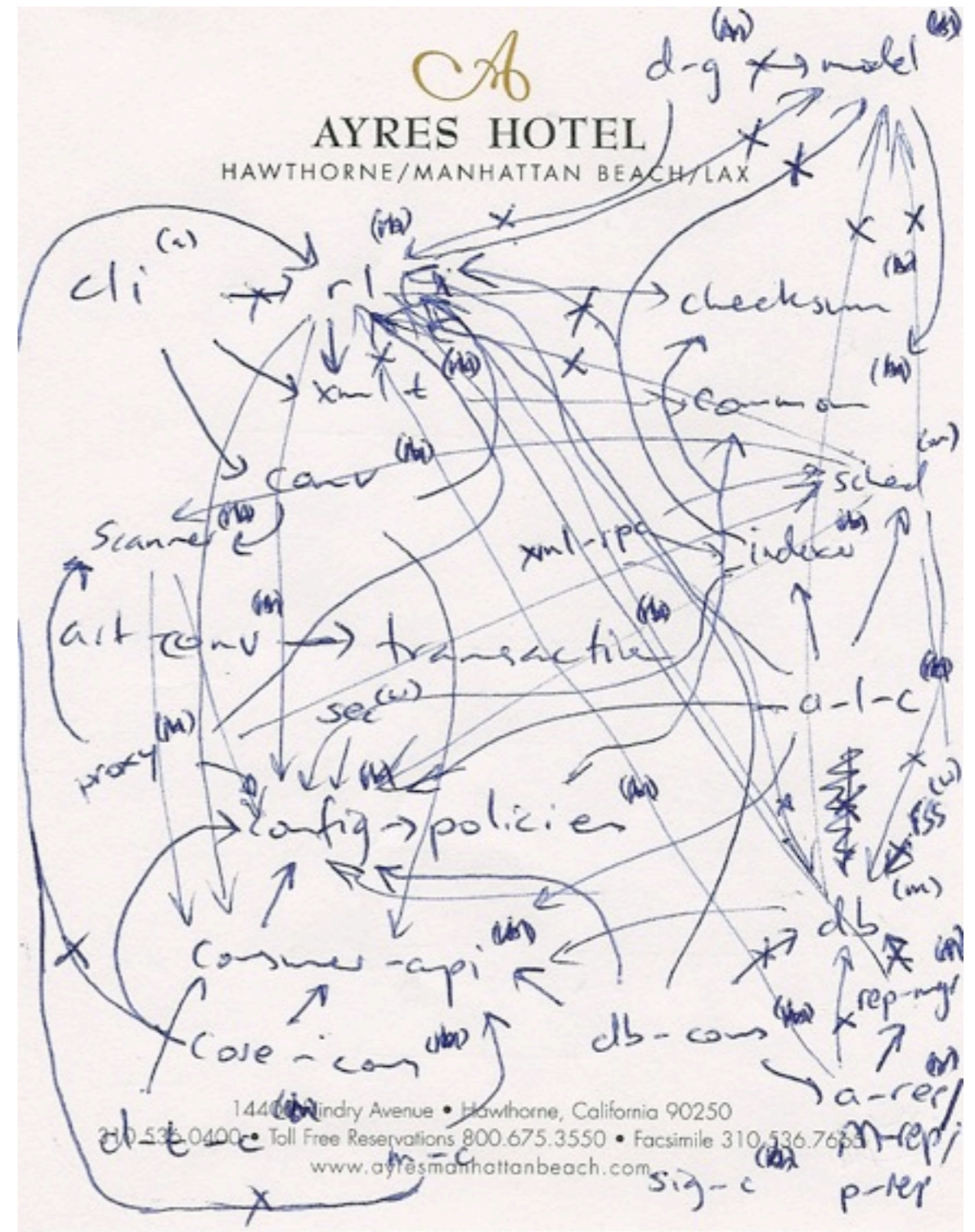
- Hierarchy familiar for XML, DOM, filesystems
- How much structure is known in advance?
- What type of queries are needed?
  - hierarchy good at locating content, but not based on joined data
- Databases are not as good at transitive retrieval, or navigation / traversal of data
  
- Archiva is hierarchical
  - filesystem-like structure
  - POM inheritance & dependency relationships

# Archiva: First Steps

- Reviewed the architecture

# Archiva: First Steps

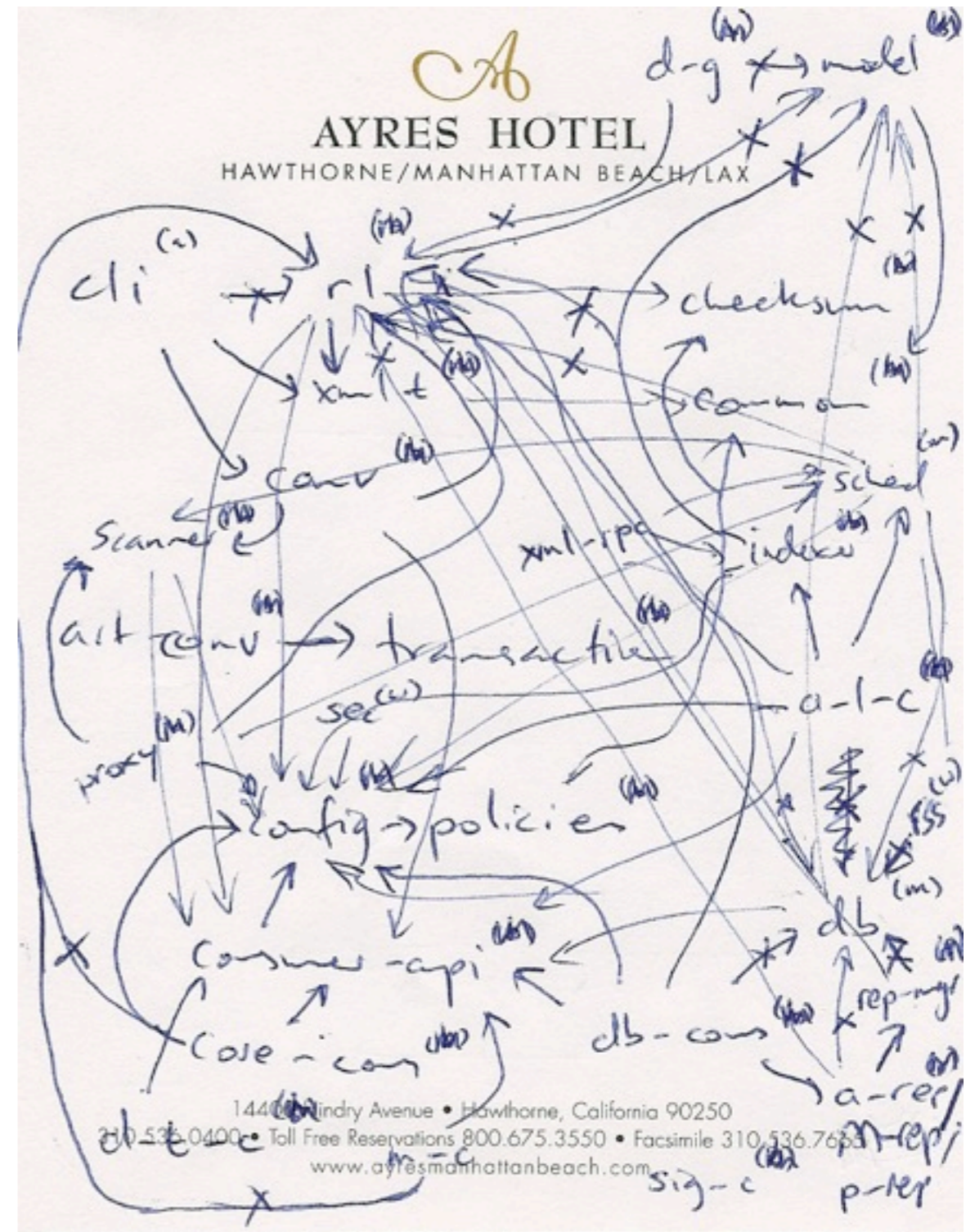
- Reviewed the architecture



# Archiva: First Steps

- Reviewed the architecture

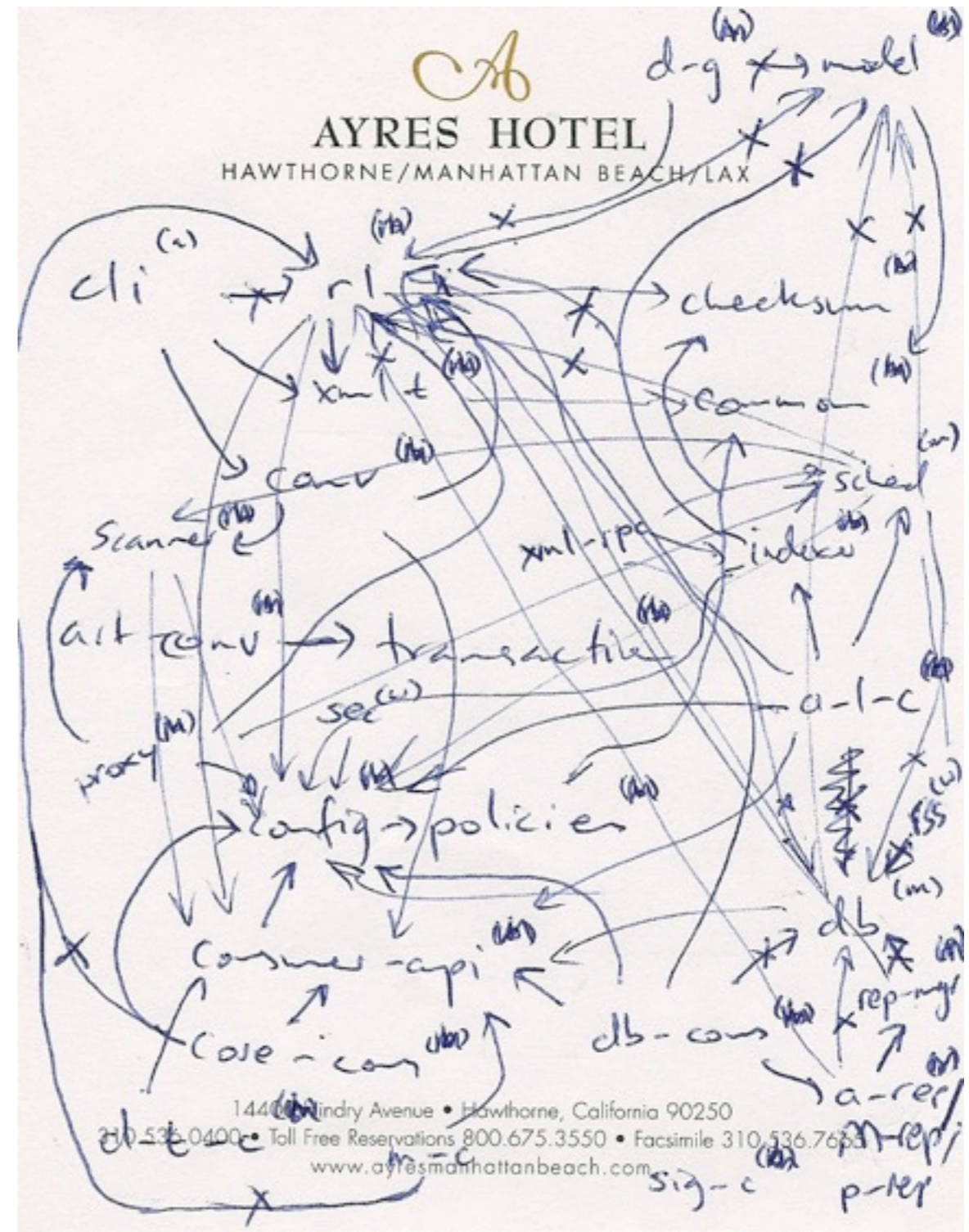
YIKES!





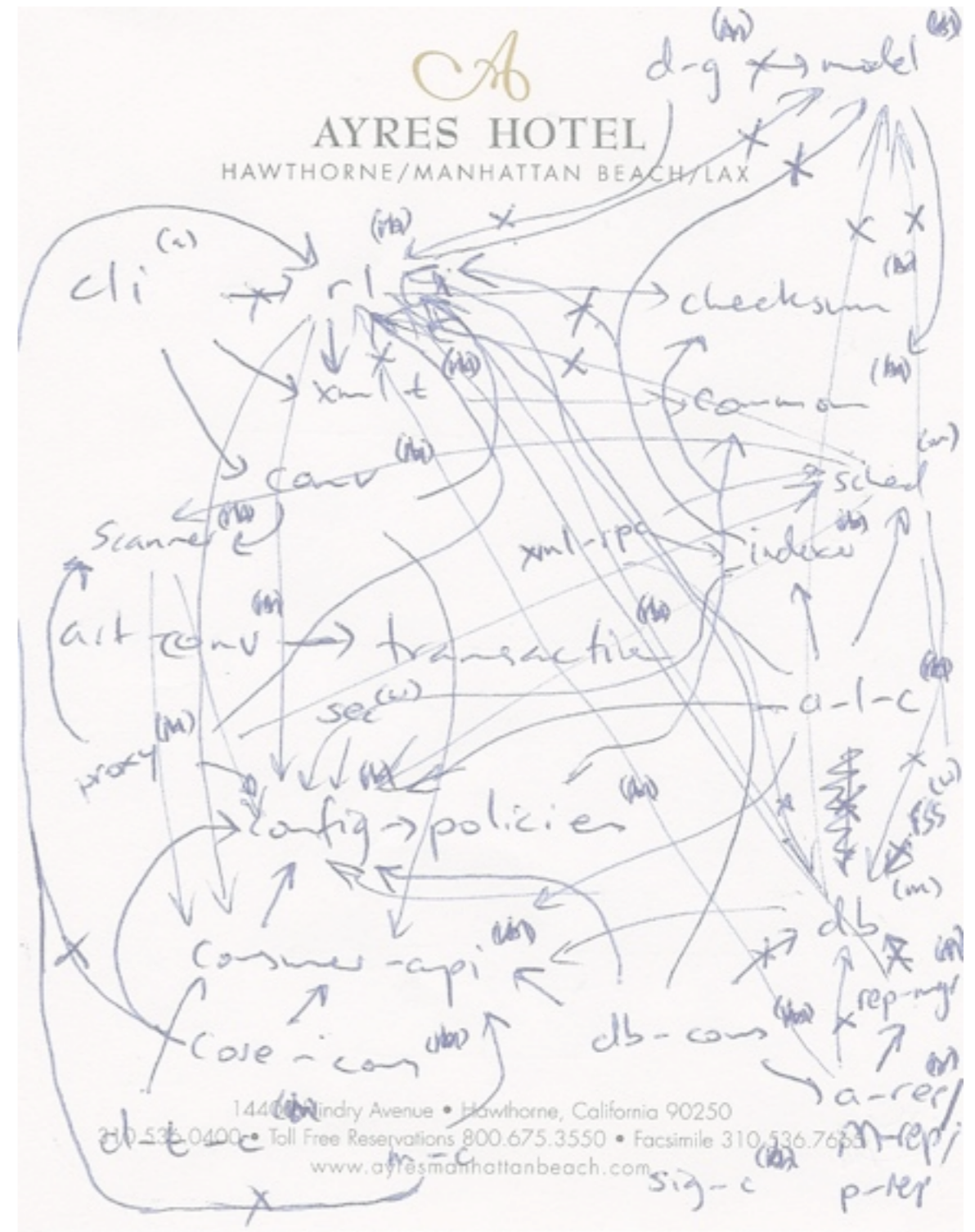
# Target Restructuring

- Technically unrelated to the effort at hand
- Trying to remove the database showed how pervasive it was
  - model (JDO annotated classes) and database module were everywhere
  - repository-layer was the culprit
  - other modules grew up out of what was available
- Started to build the right abstraction
  - directly replace some uses
  - others were redirected via the old code



# Target Restructuring

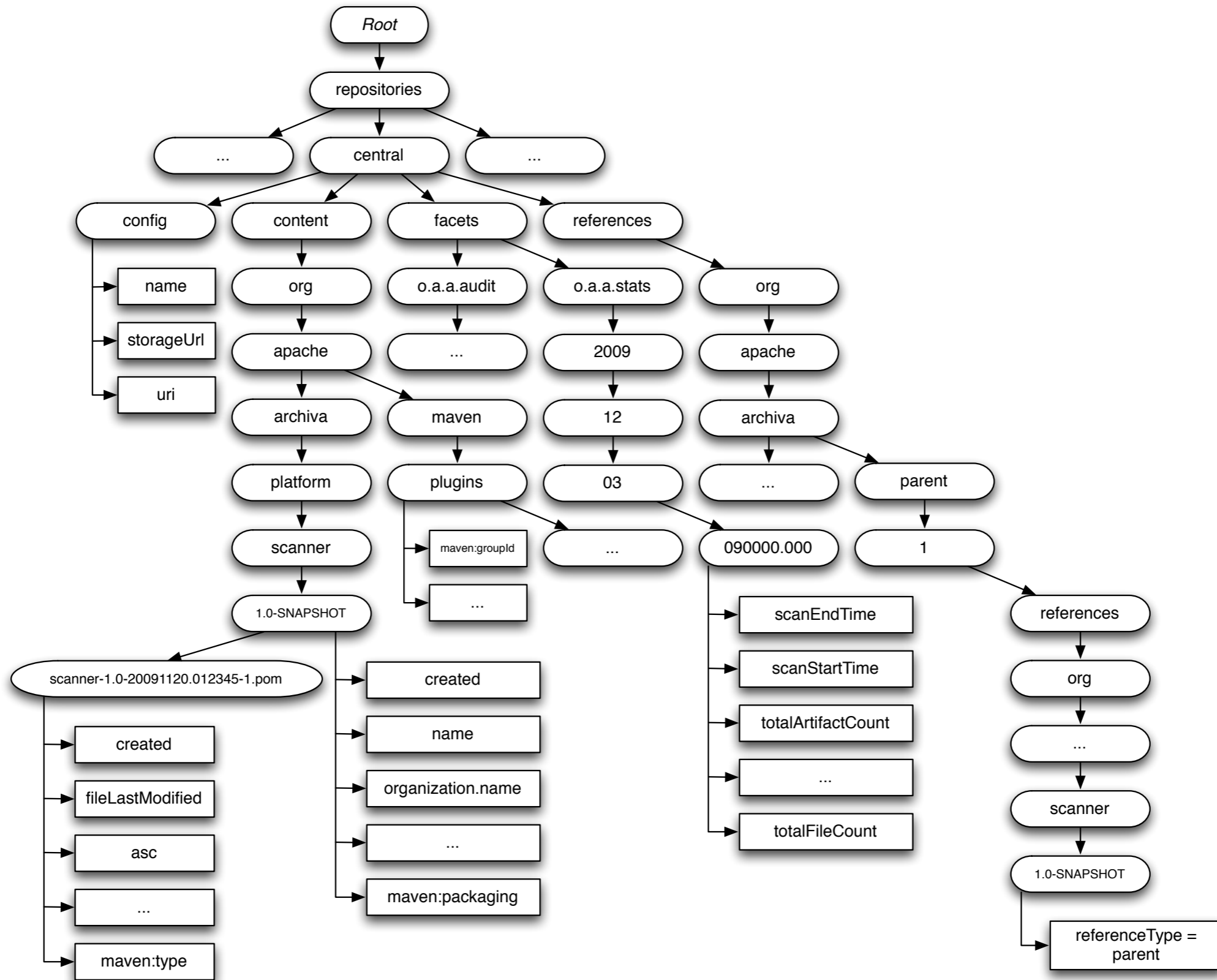
- Technically unrelated to the effort at hand
- Trying to remove the database showed how pervasive it was
  - model (JDO annotated classes) and database module were everywhere
  - repository-layer was the culprit
  - other modules grew up out of what was available
- Started to build the right abstraction
  - directly replace some uses
  - others were redirected via the old code







# Design a Content Model

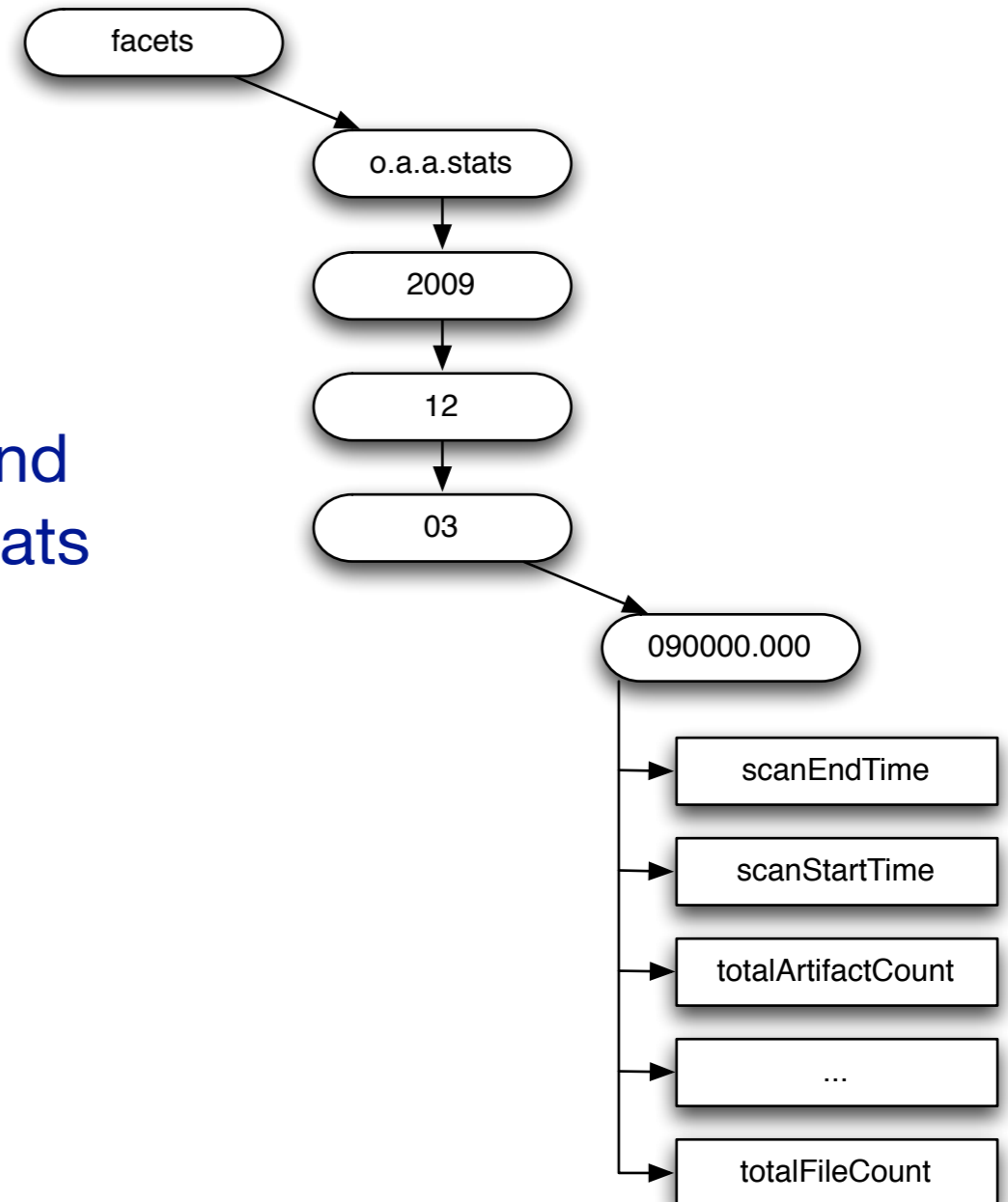


# Hierarchy

- Natural benefits for us due to the repository structure
- Can traverse to a part of a group without dealing with substrings
- Layout not identical to a Maven repository, and will be translated from a variety of input formats
- Unstructured data lends itself well to plugins and arbitrary metadata

# Hierarchy

- Natural benefits for us due to the repository structure
- Can traverse to a part of a group without dealing with substrings
- Layout not identical to a Maven repository, and will be translated from a variety of input formats
- Unstructured data lends itself well to plugins and arbitrary metadata



# David's Model

- Helpful reference
  - <http://wiki.apache.org/jackrabbit/DavidsModel>
- Rule #1: Data First, Structure Later. Maybe.
- Rule #2: Drive the content hierarchy, don't let it happen.
- Rule #3: Workspaces are for clone(), merge() and update().
- Rule #4: Beware of Same Name Siblings.
- Rule #5: References considered harmful.
- Rule #6: Files are Files are Files.
- Rule #7: ID's are evil.

# Faceted Metadata

- Designed faceted metadata model and corresponding API
- Assuming a hierarchical content model, though not yet using JCR
- For most compatibility with existing code, fully mapped object model

# Repository API

- Content access mechanism
- Individual coordinate paths passed in (group, artifact, version)
  - no need to construct the strings, avoids layout being spread out
- Resolvers to fill in metadata or obtain artifacts
  - layered access
  - track completeness
  - proxying remotely
- Resolver is a generally useful pattern if you need to load the data from an external source on the fly
- Metadata vs. Storage

# Metadata Persistence

- To get it working, a simple hand-rolled file-based implementation
- With everything working, now saw what we could achieve with JCR

# JCR - Java Content Repository

- Using Jackrabbit
- Very simple translation to the JCR API
- Initial memory usage is much higher
- Performance was still better than others
- Switched to file-based persistence of the content repository
  - `<PersistenceManager class = "org.apache.jackrabbit.core.persistence.bundle.BundleFsPersistenceManager"/>`
  - Yet to be proven at scale in Archiva, potentially not as performant
- Not using OCM (Object Content Mapping - similar to ORM but for content repositories)



# Filling Metadata

```
ProjectVersionMetadata versionMetadata =
    new ProjectVersionMetadata();

try
{
    Node root = session.getRootNode();

    Node node = root.getNode(
        "repositories/" + repositoryId + "/content/" +
        namespace + "/" + projectId + "/" + projectVersion );

    versionMetadata.setId( projectVersion );
    versionMetadata.setName(
        node.hasProperty( "name" ) ?
            node.getProperty( "name" ) :
            null );
    ...
}
```

# Adding Metadata

```
try
{
    Node root = session.getRootNode();

    Node node = root.getNode(
        "repositories/" + repositoryId + "/content/" +
        namespace + "/" + projectId );

    Node versionNode = node.addNode(
        versionMetadata.getId() );

    versionNode.setProperty( "name",
        versionMetadata.getName() );
    versionNode.setProperty( "description",
        versionMetadata.getDescription() );
    ...
}
```

# Querying

```
SELECT * FROM [archiva:artifact] AS artifact WHERE  
ISDESCENDANTNODE(artifact, '/repositories/central/content/')  
AND ([sha1] = $checksum OR [md5] = $checksum)
```

```
Query query =  
    session.getWorkspace().getQueryManager().createQuery( q,  
        Query.JCR_SQL2 );
```

```
ValueFactory valueFactory = session.getValueFactory();  
query.bindValue( "checksum", valueFactory.createValue( checksum ) );
```

```
QueryResult result = query.execute();  
artifacts = new ArrayList<ArtifactMetadata>();  
for ( Node n : JcrUtils.getNodes( result ) ) {  
    artifacts.add( getArtifactFromNode( repositoryId, n ) );  
}
```

# What about Modularisation?

- Modularisation was key to our changes
- Can be done no matter what technologies you use
- Planning for OSGi, but other priorities have come up on the way

# Where the Changes Helped

- Scanning vs. On-demand
- Dependency structure more performant and reliable
- Database was removed
  - whole class of exceptions just disappeared
  - previously unreliable operations like reverse dependency tree fixed
  - memory usage reduced
  - configuration simplified
- Metadata is more extensible
  - generic metadata plugin
  - new plugins contribute metadata without changing code or schema

# Challenges

- Initially query by artifact properties
  - e.g. how to find an artifact with a given checksum
  - starting to use JCR 2.0 queries which are being much more effective
- Correctly configuring Jackrabbit

# Opportunities

- Exposing JCR API directly to Archiva plugins
- Integration of existing WebDAV access
- Security access directly integrated into JCR
- JCR event model
- JCR version control
- OSGi
- Sling
- General design - lots more to do!

# Tips

- Review whether data is hierarchical or structure derived from the data
- Centralise access, but don't overdo the abstraction
- Align content model to natural usage
  - try not to deal with constructing and parsing paths
  - deal with content directly rather than translating to objects
- Huge value in having automated unit tests to keep it working as you make significant changes



# Archiva: Help Wanted

- Looking for developers to get involved
- [dev@archiva.apache.org](mailto:dev@archiva.apache.org)

# Thanks!

- <http://archiva.apache.org/>
- <http://archiva.apache.org/ref/1.4-M1/>
- <http://wiki.apache.org/jackrabbit/DavidsModel>
- <http://www.scribd.com/doc/11163161/JCR-or-RDBMS-why-when-how>
- <http://wiki.apache.org/jackrabbit/JcrLinks>
- <http://jackrabbit.apache.org/>
- <http://maven.apache.org/>