

# LUCENE AND BEYOND

## Core Storage Extension In OpenSearch



# LUCENE AND BEYOND

## Core Storage Extension In OpenSearch



**SAMUEL (SAM) HERMAN**

Sr Principal Engineer - Oracle



# WHY CHANGING STORAGE FORMAT?

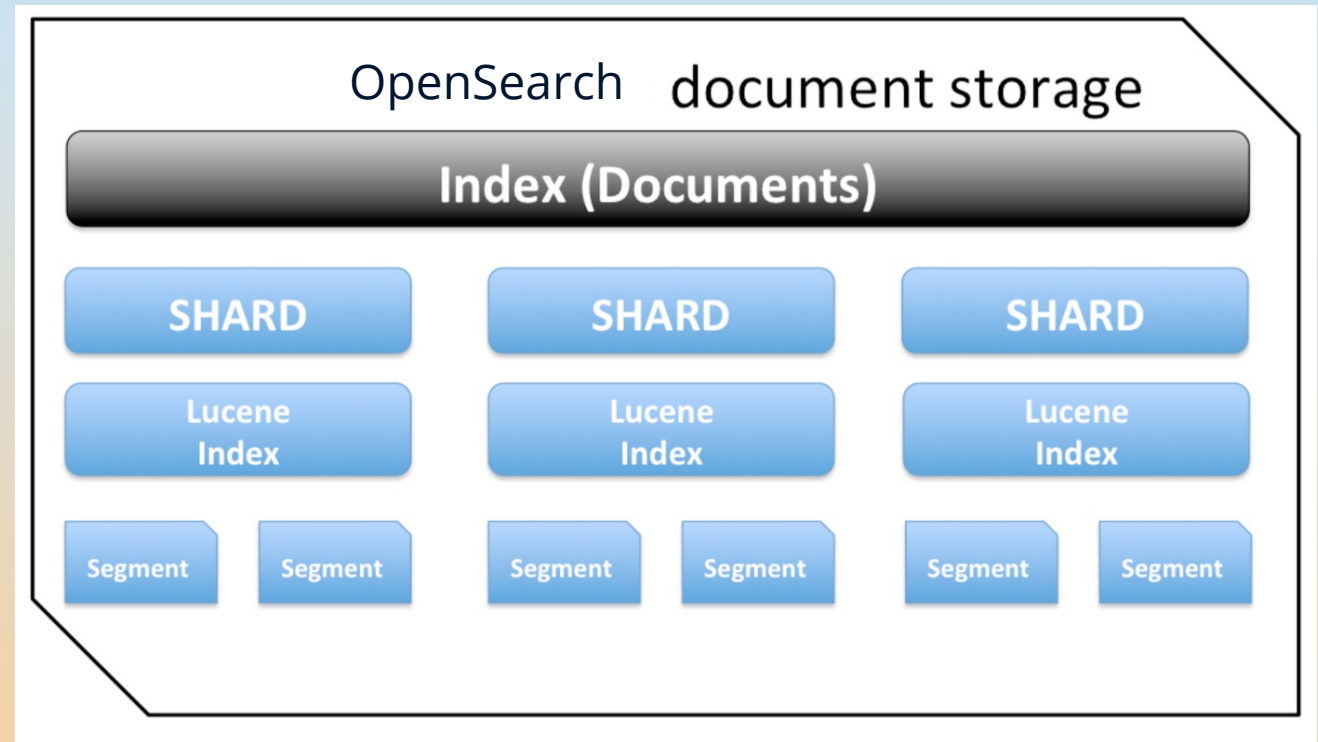
- Easier compatibility with other data warehousing solutions
- Extend functionality
- Better performance

# IMPORTANT DATA ABSTRACTIONS IN OPENSEARCH

Mapper Extension - Adding new field/data type

Engine - Database operations

Lucene Codec - Primitive types properties (e.g. Encoding/Format)



# ADDING NEW FIELD/MAPPING TYPE IN OPENSEARCH

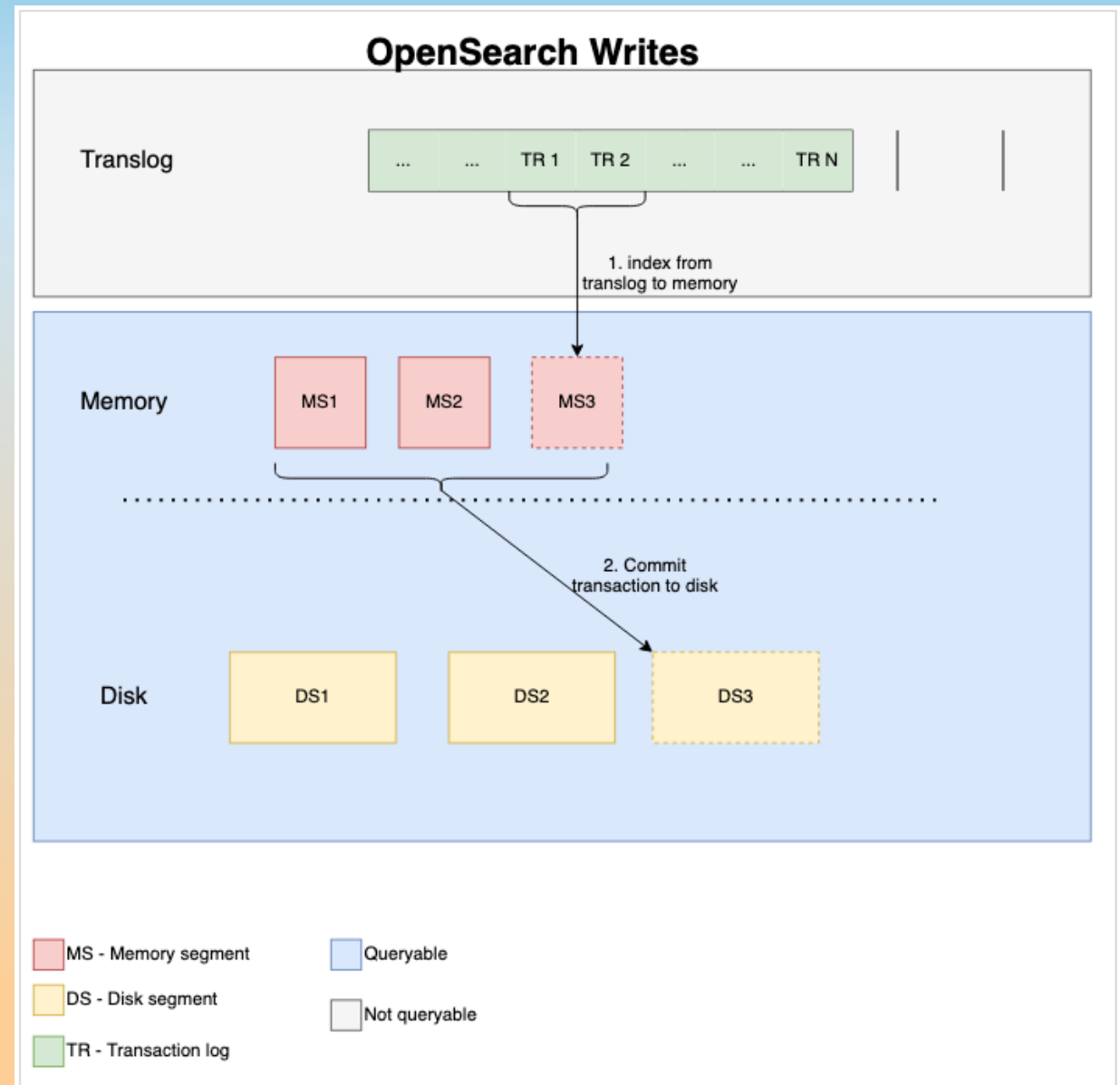
- IP, date, geo point ranges
- New field types
- Mappers and parsers
- Built upon more primitive types (e.g. Lucene)

# ENGINE

Facilitates all  
index/read/search  
operations on a SHARD level

Transactions/merges/files  
tracking/commits etc..

Entirely coupled with Lucene



# LUCENE As AN EXTENDIBLE FRAMEWORK

- Directory
- Index/search APIs
- Commits/Transactions/Segments management
- Fields
- Codec

# LUCENE – DOCUMENTS AND FIELDS

- Document is a collection of fields and values
- Field types: Text, String, long, int, double, sortedField
- Indexable/Non indexable
- Stored/Not stored
- DocValuesField/StoredField



# COMPONENTS OF A LUCENE FIELD

- Posting lists
- Doc Values
- Stored Fields
- Points
- Knn vectors

# LUCENE CODEC

How each of the field formats are implemented:

- KnnVector - implemented as a dense vector
- Point – BKD tree

# DocVALUES vs STORED FIELDS

Row storage

AAPL	2014-01-06	543.93	...	...
AAPL	2014-01-07	540.04	...	...
AAPL	2014-01-07	540.04	...	...
AAPL	2014-01-07	540.04	...	...
AAPL	2014-01-07	540.04	...	...

Column storage

AAPL	2014-01-06	543.93	...	...
AAPL	2014-01-07	540.04	...	...
GOOG	2014-01-06	558.10	...	...
GOOG	2014-01-07	568.86	...	...
GOOG	2014-01-08	570.04	...	...

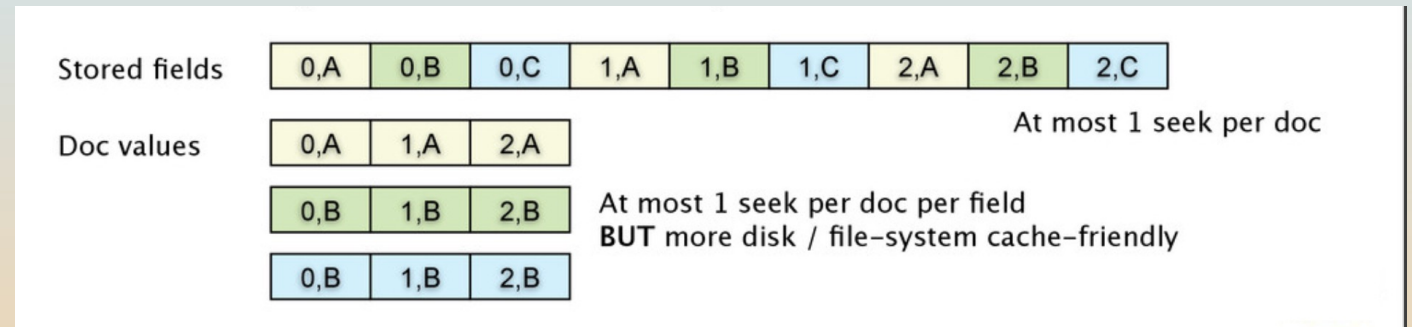
# DocVALUES Vs STORED FIELDS

## Stored fields

- get many field values for a single doc
- LZ4 compressed (default)

## Doc values

- Get few field values for many docs, good for faceting
- Various compression techniques optimized for each doc value type and also depending on the statistical properties of the data set in each block



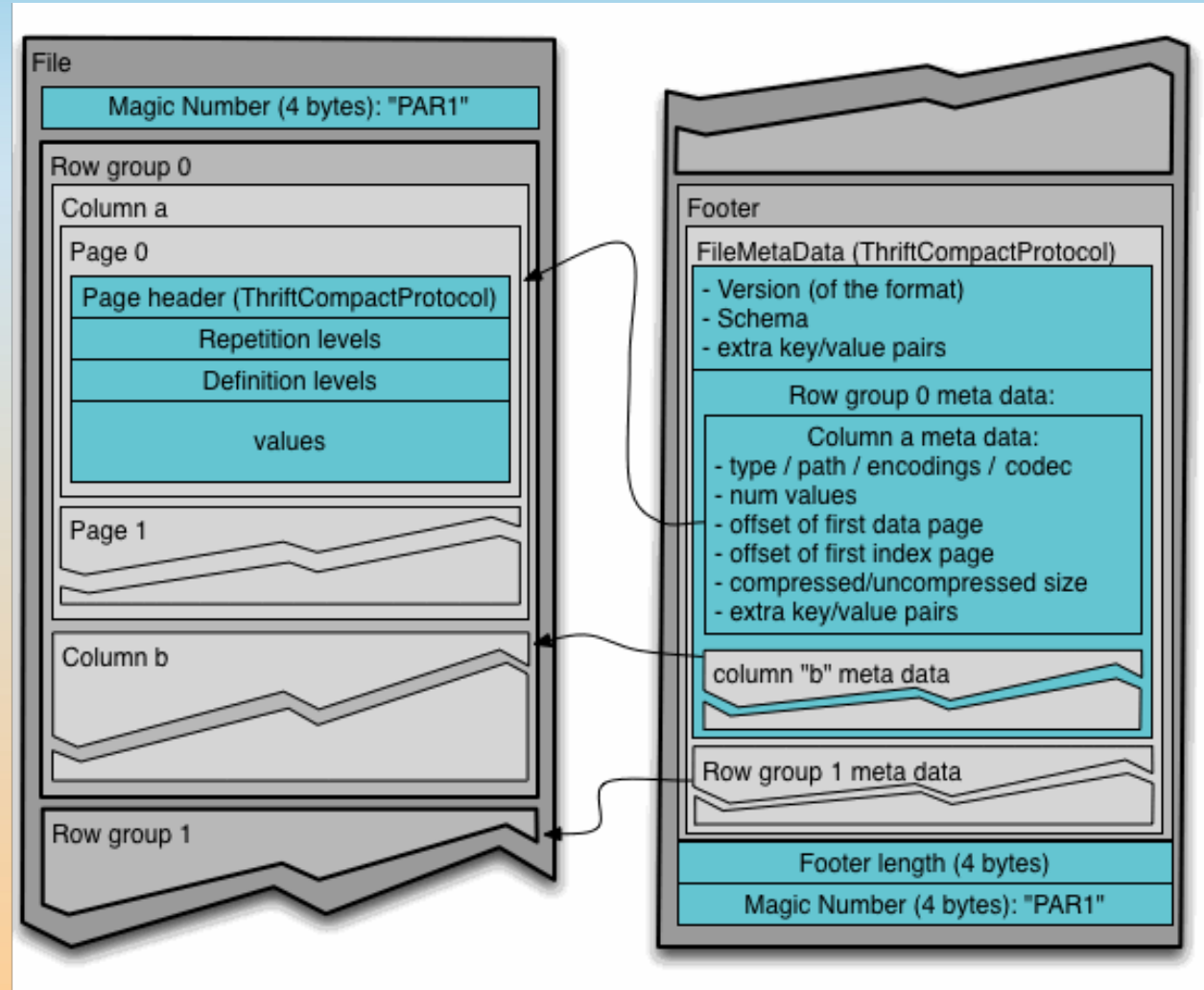
# WHAT ABOUT COMPLEX FORMATS?

Where should we fit more complex formats? Parquet, Avro etc..

# PARQUET

```

4-byte magic number "PAR1"
<Column 1 Chunk 1 + Column Metadata>
<Column 2 Chunk 1 + Column Metadata>
...
<Column N Chunk 1 + Column Metadata>
<Column 1 Chunk 2 + Column Metadata>
<Column 2 Chunk 2 + Column Metadata>
...
<Column N Chunk 2 + Column Metadata>
...
<Column 1 Chunk M + Column Metadata>
<Column 2 Chunk M + Column Metadata>
...
<Column N Chunk M + Column Metadata>
File Metadata
4-byte length in bytes of file metadata (little endian)
4-byte magic number "PAR1"
  
```



# IMPLEMENT As DocVALUES?

- Consumer / Producer interface
- Different consumer for every field
- Optimized for throughput, writes all the doc values of a field in one go
- Limiting us in preserving row group structure

# IMPLEMENT AS DocVALUES WITH BINARY FIELDS?

- Treat as a single stored binary doc value that passes a record
- Client behavior change – awareness of data format
- OpenSearch with Lucene can be used to enhance the interpretation of the field



# IMPLEMENT As STORED FIELDS?

- Preserves row group structure
- Easier to read outside of the cluster
- Double buffering
- Misleading the abstraction implementation

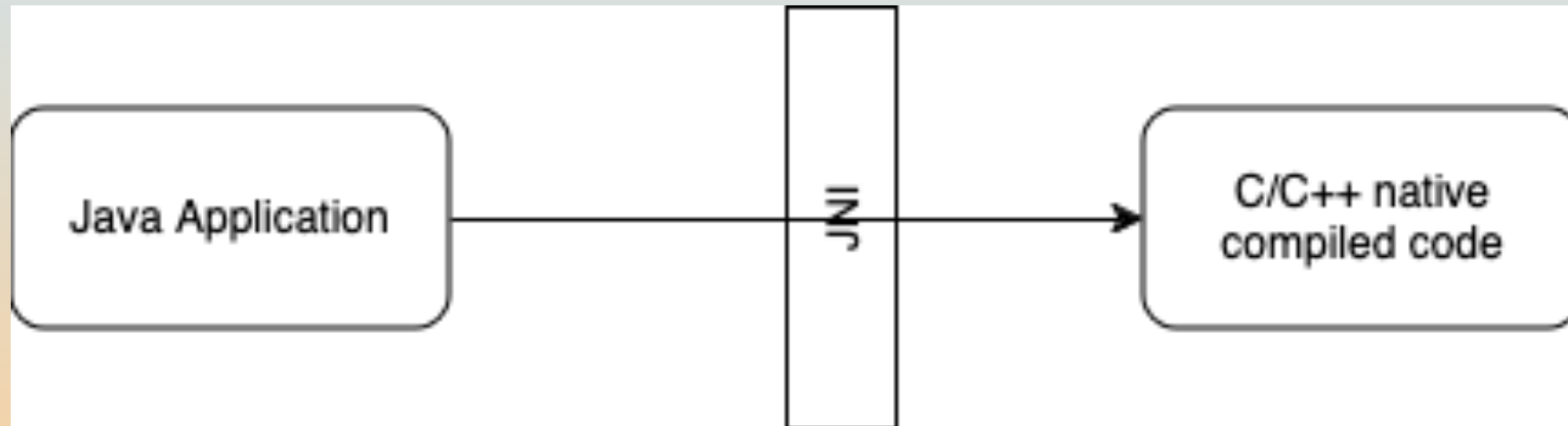
# EXTEND THE LUCENE INDEXINGCHAIN OR OTHER INTERFACES

- Gap in Lucene interfaces
- IndexingChain can be extended
- Field can be extended

# LIMITATIONS

- Lucene - Java
- OpenSearch – Java + JarHell

# POSSIBLE SOLUTIONS



# EXTENSIONS

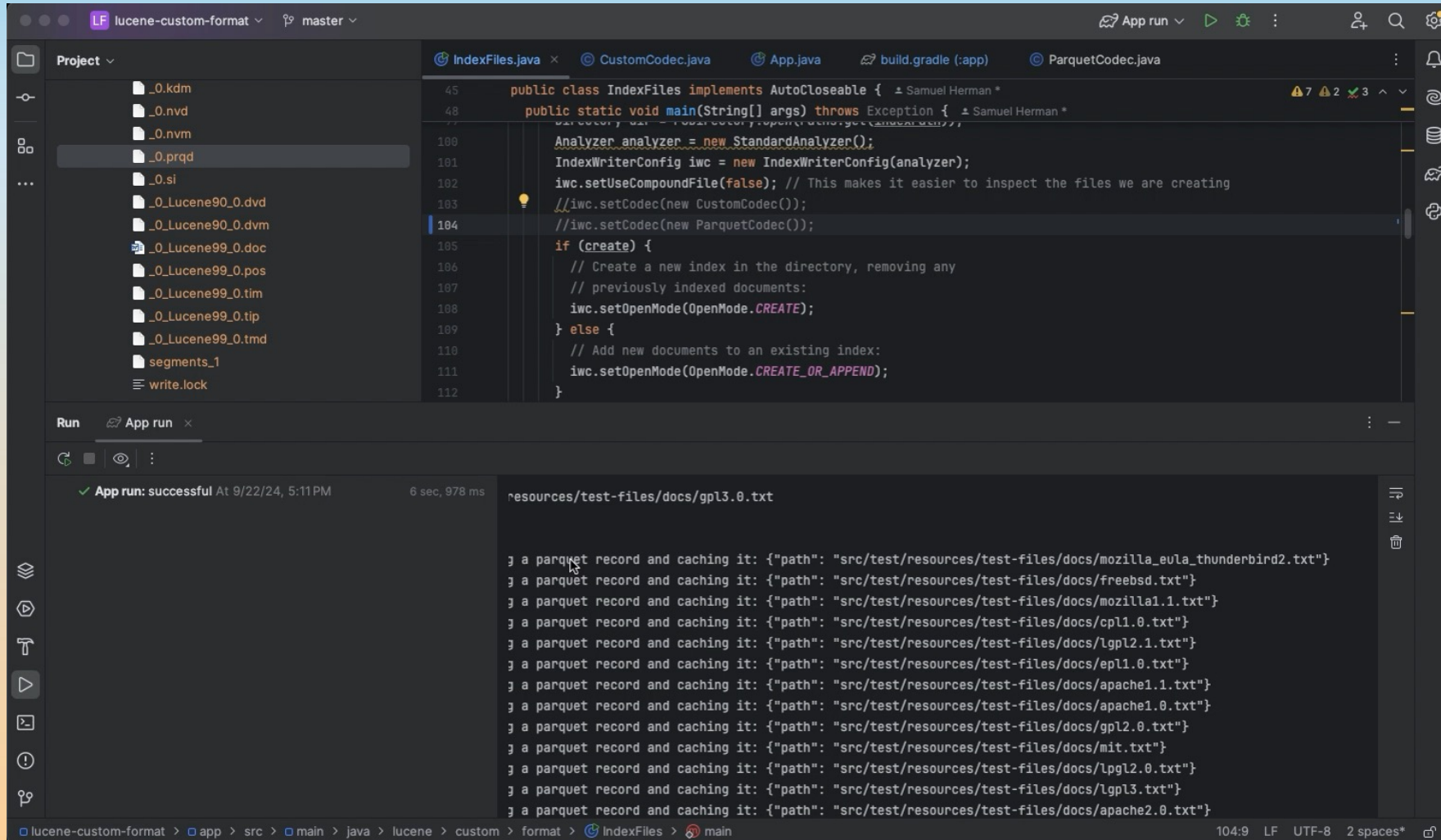
- Experimental
- Not working with local storage
- Launch new process and make sure it's local
- Route all requests based on locality

<https://github.com/opensearch-project/OpenSearch/issues/2447>

# LOCAL EXTERNAL WRITER

- Can be extended for non Java formats
- Solves the issues of locality
- Solves the issues of dependency collision

<https://github.com/opensearch-project/OpenSearch/issues/13668>



The screenshot shows an IDE window for a project named 'lucene-custom-format'. The left sidebar displays a file tree with various test files, including '\_0.kdm', '\_0.nvd', '\_0.nvm', '\_0.prqd', '\_0.si', and several '\_0\_Lucene90' and '\_0\_Lucene99' files. The main editor area shows the 'IndexFiles.java' file, which implements the 'AutoCloseable' interface. The code defines a 'main' method that takes an array of arguments and throws an exception. It initializes an 'Analyzer' and an 'IndexWriterConfig' (iwc), setting 'iwc.setUseCompoundFile(false)' and 'iwc.setOpenMode(OpenMode.CREATE)'. The code then iterates over the arguments, creating or appending documents to the index. The output window at the bottom shows the successful execution of the application, displaying the path 'resources/test-files/docs/gpl3.0.txt' and a list of documents being indexed, such as 'mozilla\_eula\_thunderbird2.txt', 'freebsd.txt', 'mozilla1.1.txt', 'cpl1.0.txt', 'lpgl2.1.txt', 'epl1.0.txt', 'apache1.1.txt', 'apache1.0.txt', 'gpl2.0.txt', 'mit.txt', 'lpgl2.0.txt', 'lpgl3.txt', and 'apache2.0.txt'.

```
public class IndexFiles implements AutoCloseable {
    public static void main(String[] args) throws Exception {
        Analyzer analyzer = new StandardAnalyzer();
        IndexWriterConfig iwc = new IndexWriterConfig(analyzer);
        iwc.setUseCompoundFile(false); // This makes it easier to inspect the files we are creating
        iwc.setCodec(new CustomCodec());
        iwc.setOpenMode(OpenMode.CREATE);
        if (create) {
            // Create a new index in the directory, removing any
            // previously indexed documents:
            iwc.setOpenMode(OpenMode.CREATE);
        } else {
            // Add new documents to an existing index:
            iwc.setOpenMode(OpenMode.CREATE_OR_APPEND);
        }
    }
}
```

App run: successful At 9/22/24, 5:11 PM 6 sec, 978 ms

resources/test-files/docs/gpl3.0.txt

g a parquet record and caching it: {"path": "src/test/resources/test-files/docs/mozilla\_eula\_thunderbird2.txt"}

g a parquet record and caching it: {"path": "src/test/resources/test-files/docs/freebsd.txt"}

g a parquet record and caching it: {"path": "src/test/resources/test-files/docs/mozilla1.1.txt"}

g a parquet record and caching it: {"path": "src/test/resources/test-files/docs/cpl1.0.txt"}

g a parquet record and caching it: {"path": "src/test/resources/test-files/docs/lpgl2.1.txt"}

g a parquet record and caching it: {"path": "src/test/resources/test-files/docs/epl1.0.txt"}

g a parquet record and caching it: {"path": "src/test/resources/test-files/docs/apache1.1.txt"}

g a parquet record and caching it: {"path": "src/test/resources/test-files/docs/apache1.0.txt"}

g a parquet record and caching it: {"path": "src/test/resources/test-files/docs/gpl2.0.txt"}

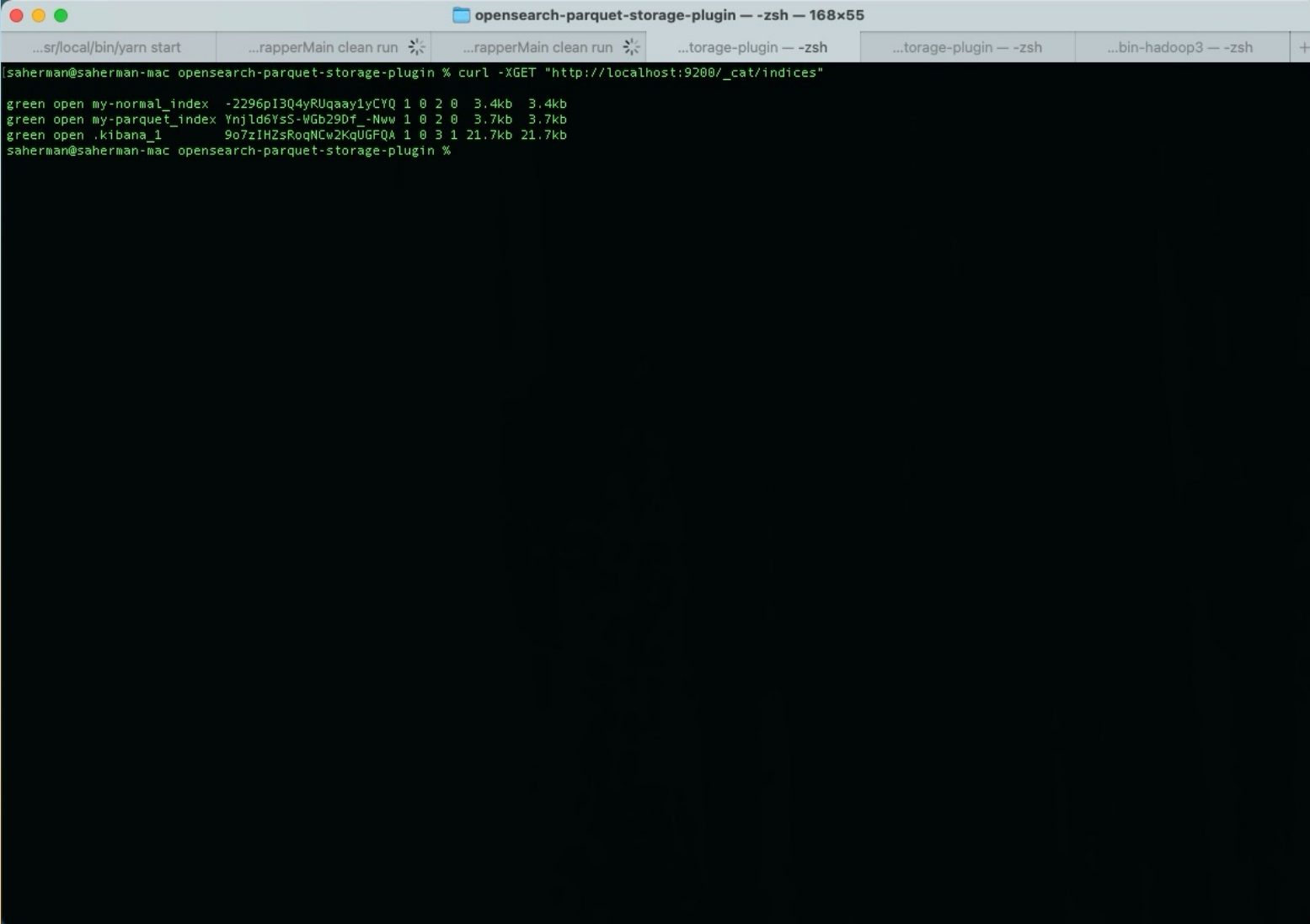
g a parquet record and caching it: {"path": "src/test/resources/test-files/docs/mit.txt"}

g a parquet record and caching it: {"path": "src/test/resources/test-files/docs/lpgl2.0.txt"}

g a parquet record and caching it: {"path": "src/test/resources/test-files/docs/lpgl3.txt"}

g a parquet record and caching it: {"path": "src/test/resources/test-files/docs/apache2.0.txt"}

# LUCENE CUSTOM CODEC DEMO (AND PARQUET)



A terminal window titled "opensearch-parquet-storage-plugin -- zsh -- 168x55" is shown. The terminal displays the output of a curl command to the OpenSearch \_cat/indices endpoint. The output shows three indices: "my-normal\_index", "my-parquet\_index", and ".kibana\_1". Each index is in a "green" status and "open". The "my-parquet\_index" is highlighted with a green background. The terminal also shows the prompt "saherman@saherman-mac opensearch-parquet-storage-plugin %".

```
saherman@saherman-mac opensearch-parquet-storage-plugin % curl -XGET "http://localhost:9200/_cat/indices"

green open my-normal_index -2296pI3Q4yRUqaay1yCYQ 1 0 2 0 3.4kb 3.4kb
green open my-parquet_index Ynjld6YsS-WGb29Df_-Nww 1 0 2 0 3.7kb 3.7kb
green open .kibana_1 9o7zIHZsRoqNCw2KqUGFQA 1 0 3 1 21.7kb 21.7kb
saherman@saherman-mac opensearch-parquet-storage-plugin %
```

# OPENSEARCH ON PARQUET!



# QUESTIONS?