# Apache Lucene:
# From Text Indexing to Artificial Intelligence

**Lucian Precup**

2024-04-23

#HaystackConf

@lucianprecup

@a2lean

1

# Poll
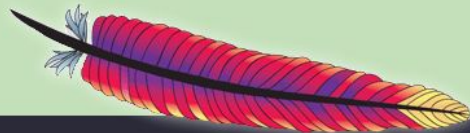
- Apache Lucene?
- Apache Solr? Elasticsearch? OpenSearch?
- Other search technologies?

2024-04-23          #HaystackConf          @a2lean

# Apache Lucene - 22.5 years and counting

# Apache Lucene - 22.5 years and counting



**Doug Cutting** @cutting

Lucene's FuzzyQuery is 100 times faster!
blog.mikemccandless.com/2011/03/lucene...

9:58 PM · Mar 24, 2011

**Uwe 🇺🇦 Schindler** 👨‍🦰🇩🇪🪶🪶🪶 @thetaph1

Mike McCandless talks about #Apache #Lucene that helps to squash #Java #JVM bugs: elastic.co/blog/lucene-jv...

9:45 PM · Jul 17, 2015 from Bremen, Germany

**Adrien Grand** @jpountz

I ran some benchmarks between Lucene 9.7 and 9.8 (soon to be released), as well as with recursive graph partitioning enabled (-bp): jpountz.github.io/lucene-9.7-vs-... There's a nice speedup on 9.8 alone, and then recursive graph bisection gives another great speedup.

11:46 PM · Sep 13, 2023 · **3,107** Views

**Uwe 🇺🇦 Schindler** 👨‍🦰🇩🇪🪶🪶🪶 @thetaph1

#Apache #Lucene can much faster execute kNN vector queries by calculating dot products / cosine distances using SIMD instructions on AVX2 (x86) and NEON (ARM). It will only work with #Java20 on coming Lucene 9.7 with "--add-modules jdk.incubator.vector":

apache/lucene
**#12311 Integrate the Incubating Panama Vector API**

💬 165 comments   🔲 62 reviews   ± 16 files   **+1025 -179** ■■■■

ChrisHegarty · May 18, 2023  ⦿ 59 commits

github.com
Integrate the Incubating Panama Vector API by ChrisHegarty · Pull Request #...
Leverage accelerated vector hardware instructions in Vector Search. Lucene already has a mechanism that enables the use of non-final JDK APIs, currenti...

7:26 PM · May 26, 2023 · **23.1K** Views

# Berlin Buzzwords - The Search, Store and Scale Conference



Berlin Buzzwords 2015: Uwe Schindler -
Apache Lucene 5 - ...
598 views · Jun 3, 2015
YouTube › Plain Schwarz



Berlin Buzzwords 2016: Uwe Schindler -
Apache Lucene 6: W...
704 views · Jun 12, 2016
YouTube › Plain Schwarz



Berlin Buzzwords 2017: Uwe Schindler -
Apache Lucene 7 - ...
741 views · Jun 15, 2017
YouTube › Plain Schwarz



Berlin Buzzwords 2018: Uwe Schindler –
Apache Lucene an...
614 views · Jun 13, 2018
YouTube › Plain Schwarz



Berlin Buzzwords 2019: Uwe Schindler –
10 (funny) years of ...
518 views · Jun 20, 2019
YouTube › Plain Schwarz

# 2031 at Berlin Buzzwords



**HAYSTACK**     2024-04-23     #HaystackConf     @a2lean

adelean
EXTRACT TRANSFORM SEARCH

# Thanks to Apache Lucene

Apache Nutch – provides web crawling and HTML parsing

Apache Solr – an enterprise search server

Elasticsearch – an enterprise search server released in 2010

MongoDB Atlas Search – a cloud-native enterprise search application based on MongoDB and Apache Lucene

OpenSearch – an open source enterprise search server based on a fork of Elasticsearch 7

Adelean a2 - an e-commerce and community search server

# Who are we?



Adelean

- Experts in **search** technologies
- Integrators of **Elasticsearch, OpenSearch** and **Solr**
- **Consulting** and **Training** providers
- Developers of **a2** E-Commerce and Enterprise Search solution
- Developers of **all.site** - your **Collaborative** Search Engine

adelean
EXTRACT TRANSFORM SEARCH

# Inverted index

D1: "All economic sectors, from commerce to the automobile sector to the energy sector, are affected."

D2: "Registrations of new electric cars in France increased by 0.90% in raw data in July."

D3: "They rely on technologies that France masters on an industrial and economic level."

adelean
EXTRACT TRANSFORM SEARCH

# Inverted index

« All economic sectors, from commerce to the automobile sector to the energy sector, are affected. »

↓ Extract all words

All, economic, sectors, from, commerce, to, the, automobile, sector, to, the, energy, sector, are, affected

↓ Sort unique words

All, affected, are, automobile, commerce, economic, energy, from, sector, sector, sectors, the, to

↓ List all documents containing each word

HAYSTACK          2024-04-23          #HaystackConf          @a2lean

adelean
EXTRACT TRANSFORM SEARCH

# Inverted index

| | D1 | D2 | D3 |
|---|---|---|---|
| 0.90% | | X | |
| affected | X | | |
| All | X | | |
| an | | | X |
| and | | | X |
| are | X | | |
| automobile | X | | |
| by | | X | |
| cars | | X | |
| commerce | X | | |
| data | | X | |
| economic | X | | X |
| electric | | X | |
| energy | X | | |
| France | | X | X |
| from | X | | |
| in | | X | |
| increased | | X | |

| | D1 | D2 | D3 |
|---|---|---|---|
| industrial | | | X |
| July | | X | |
| level | | | X |
| masters | | | X |
| new | | X | |
| of | | X | |
| on | | | X |
| raw | | X | |
| Registrations | | X | |
| rely | | | X |
| sector | X | | |
| sectors | X | | |
| technologies | | | X |
| that | | | X |
| the | X | | |
| They | | | X |
| to | X | | |

adelean
EXTRACT TRANSFORM SEARCH

# Inverted index

| | D1 | D2 | D3 |
|---|---|---|---|
| 0.90% | | X | |
| affected | X | | |
| All | X | | |
| an | | | X |
| and | | | X |
| are | X | | |
| automobile | X | | |
| by | | X | |
| cars | | X | |
| commerce | X | | |
| data | | X | |
| economic | X | | X |
| electric | | X | |
| energy | X | | |
| France | | X | X |
| from | X | | |
| in | | X | |
| increased | | X | |

| | D1 | D2 | D3 |
|---|---|---|---|
| industrial | | | X |
| July | | X | |
| level | | | X |
| masters | | | X |
| new | | X | |
| of | | X | |
| on | | | X |
| raw | | X | |
| Registrations | | X | |
| rely | | | X |
| sector | X | | |
| sectors | X | | |
| technologies | | | X |
| that | | | X |
| the | X | | |
| They | | | X |
| to | X | | |

↓ Lowercase

↓ Stemming

↓ Synonyms

↓ Stopwords

↓ …

HAYSTACK          2024-04-23          #HaystackConf          @a2lean

adelean
EXTRACT TRANSFORM SEARCH

# Inverted index

|  | D1 | D2 | D3 |
|---|---|---|---|
| 0.90% |  | X |  |
| affected | X |  |  |
| all | X |  |  |
| automobile, car | X | X |  |
| commerce | X |  |  |
| data |  | X |  |
| economic | X |  | X |
| electric |  | X |  |
| energy | X |  |  |
| france |  | X | X |
| increased |  | X |  |
| industrial |  |  | X |
| july |  | X |  |
| level |  |  | X |
| master |  |  | X |
| new |  | X |  |
| raw |  | X |  |
| registration |  | X |  |
| rely |  |  | X |
| sector | X |  |  |
| technology |  |  | X |

Search : automobile ☐ Response : D1, D2
Search : France ☐ Response : D2, D3
Search : automobile AND France ☐ Response : D2

- Simple and efficient model for Boolean calculations (intersections)

- \+ Each token also contains associated metadata

# Lucene syntax

| Example | Usage |
|---|---|
| lastname:dujardin<br>firstname:jean | Prefix the search term with the column name |
| lastname:dujardin AND firstname:jean | Combine the criteria with the keyword AND (in capital letters) |
| lastname:dujardin~ | Add a "~" at the end of a term to enable fuzzy searching. |
| lastname:dujardin~0.8 | Add a ratio (~0.x) to widen/narrow the approximate search with a difference of 1-x characters out of 10 (0.9 for 1 character out of 10, 0.8 for 2 characters out of 10, etc.) |
| postalcode:59*<br>location:*yvelines | Use the "*" (star) character for a "starts with" or "containing" search. This principle does not apply for date type columns. The character "?" (question mark) can also be used to specify that it is a single missing character. |
| location:"rueil malmaison"<br>firstname:"jean françois" | Use quotation marks "" to search for a "phrase" (two or more adjacent words) |
| firstname:"jean paul"~10 | Approximate sentences. Use ~n at the end of an expression (sentence) to indicate that there may be interspersed words (n-1 words in this case). This is particularly useful for finding first names like "JEAN MARIE PAUL", "PAUL ANDRE JEAN" or even "PAUL JEAN" when searching for "jean paul" |
| birthdate:[1949-01-01 TO 1949-12-31]<br>postalcodel:[91000 TO 95000] | Use right parentheses ("[" and "]") and the TO keyword to specify an inclusive interval |
| birthdate:[1949-01-01 TO 1950-01-01} | Use curly braces ("{" and "}") to specify an exclusive interval. You can combine inclusive and exclusive intervals ("[" and "}") |
| -firstname:"jean paul" | Use the "-" character to exclude responses. In this example we exclude those who have exactly "jean paul" in their first name |

**HAYSTACK**    2024-04-23    #HaystackConf    @a2lean

adelean
EXTRACT TRANSFORM SEARCH

# Lucene text analysis

## Index time

Input documents ▢

| Id | Name |
|----|------|
| 1 | Céline |
| 2 | Celia |

Ascii folding ▢     Celine, Celia

Lowercase ▢     celine, celia

## Query time

| Name |
|------|
| CÉLINE |

▢ Search term

CELINE     ▢ Ascii folding

celine     ▢ Lowercase

## Index

| Key | Document id |
|-----|-------------|
| celine | 1 |
| celia | 2 |

HAYSTACK          2024-04-23          #HaystackConf          @a2lean

adelean
EXTRACT TRANSFORM SEARCH

# Auto-completion with nGram

## Index time

Input documents ▯

| Id | Nom |
|----|-----|
| 1 | Céline |
| 2 | Celia |

Ascii folding ▯    Celine, Celia
Lowercase ▯    celine, celia
EdgeNGram ▯    ce cel celi celin celine
ce cel celi celia

### Index

| Key | Document id |
|-----|-------------|
| ce | 1, 2 |
| cel | 1, 2 |
| celi | 1, 2 |
| celin | 1 |
| celine | 1 |
| celia | 2 |

## Query time

| Nom |
|-----|
| Célin |

▯ Search term

Celin ▯ Ascii folding
celin ▯ Lowercase
ce cel celi celin ▯ EdgeNGram

Célin

Mlle **Céline** Dupont
Mlle **Celia** Leclerc

HAYSTACK

adelean
EXTRACT TRANSFORM SEARCH

# Different analysis chain at index and query time

## Index time

Input documents ▯

| Id | Nom |
|----|-------|
| 1  | Céline |
| 2  | Celia |

Ascii folding ▯     Celine, Celia
Lowercase ▯        celine, celia
EdgeNGram ▯     ce cel celi celin celine
                         ce cel celi celia

## Query time

| Nom |
|------|
| Célin |

▯ Search term

Celin ▯ Ascii folding
celin ▯ Lowercase

### Index

| Key | Document id |
|--------|-------------|
| ce     | 1, 2 |
| cel    | 1, 2 |
| celi   | 1, 2 |
| celin  | 1 |
| celine | 1 |
| celia  | 2 |

Celin

Mlle **Céline** Dupont

HAYSTACK

adelean
EXTRACT TRANSFORM SEARCH

# Analysis menu in the Solr Admin dashboard

# Automaton Invasion

FSA/Ts were difficult to implement but proved to be a good match for Lucene

FSA/Ts are now used in a number of places

FuzzyQuery: 100x times faster

SynonymFilter: 5x times faster, 59X less RAM

https://blog.mikemccandless.com/2011/03/lucenes-fuzzyquery-is-100-times-faster.html

https://blog.mikemccandless.com/2012/05/finite-state-automata-in-lucene.html

# Adding (real) support for numbers



From sorting by pseudo-numbers :

"00000000010" < "00000000300"☺️

to supporting numbers :

10 < 300

# Facets, facets, facets, …



Term facets

Range facets

Hierarchical facets

…

# Lucene 4.0: From FieldCache to DocValues

| weight | |
|:---:|:---|
| 5.8 | |
| 1.0 | |
| 2.7 | ← |
| 2.7 | ← |
| 4.3 | |
| 7.9 | |
| 1.0 | |
| 3.2 | |
| 4.7 | |
| 7.9 | |
| 9.0 | |

parse → convert to datatype

array per field / segment

float 32

un-invert

| term | freq | Posting list |
|:---:|:---:|:---|
| 1.0 | 1 | 1 6 |
| 2.7 | 1 | 2 3 |
| 3.2 | 1 | 7 |
| 4.3 | 1 | 4 |
| 4.7 | 1 | 8 |
| 5.8 | 1 | 0 |
| 7.9 | 1 | 5 9 |
| 9.0 | 1 | 10 |

string / byte[]

→ Excellent performance of sorting, aggregations and highlighting
→ Lucene: a columnar store

| | 100k Docs | 1M Docs | 10M Docs |
|:---|:---:|:---:|:---:|
| FieldCache | 122 ms | 348 ms | 3161 ms |
| DocValues | **7 ms** | **10 ms** | **90 ms** |

https://www.slideshare.net/lucenerevolution/willnauer-simon-doc-values-column-stride-fields-in-lucene

# Apache Lucene - a relational database?



```json
{
    "film" : {
        "id" : "183070",
        "title" : "The Artist",
        "published" : "2011-10-12",
        "genre" : ["Romance", "Drama", "Comedy"],
        "language" : ["English", "French"],
        "persons" : [
            {"person" : { "id" : "5079",  "name"
: "Michel Hazanavicius",  "role" : "director" }},
            {"person" : { "id" : "84145",  "name"
: "Jean Dujardin",  "role" : "actor" }},
            {"person" : { "id" : "24485",  "name"
: "Bérénice Bejo",  "role" : "actor" }},
            {"person" : { "id" : "4204",  "name"
: "John Goodman",  "role" : "actor" }}
        ]
    }
}
```

# Block joins and result grouping in Apache Lucene

https://www.slideshare.net/lucenerevolution/grouping-and-joining-in-lucenesolr

https://cwiki.apache.org/confluence/display/solr/FieldCollapsing

# Multi-dimensional points, coming in Apache Lucene 6.0



https://www.elastic.co/blog/lucene-points-6-0

# Lucene - a spatial database

Excellent support for Geo Shape (dots, lines, polygons)

Spatial data requesting (Geo Shape Query, Geo Bounding Box Query, Geo Distance Query, Geo Polygon Query)

This allowed the development of Kibana Maps and the user-friendly support for spatial data in Kibana

# New use cases: Monitoring, Observability, ...

# Data analysis with Kibana 🚀


Discover


Lens


Dashboards


Alerting


Maps


Canvas


Machine Learning


Graph

HAYSTACK          2024-04-23          #HaystackConf          @a2lean

adelean
EXTRACT TRANSFORM SEARCH

# NLP models and vector search: a paradigm shift

Move beyond text-matching

https://www.elastic.co/elasticsearch/vector-database

# NLP models and vector search: a paradigm shift

Move beyond text-matching

https://www.elastic.co/elasticsearch/vector-database

# Embeddings or dense semantic vector

Vectorization makes it possible to represent textual documents as numbers in a mathematical space, which facilitates their analysis and processing by computer algorithms.

Each element in this vector represents a characteristic of the document.

MAMMAL

REALISTIC

CARTOON

BIRD

https://www.elastic.co/blog/whats-new-elasticsearch-kibana-cloud-8-4-0

adelean
EXTRACT TRANSFORM SEARCH

# Vectorization in a three-dimensional vector



| Animal | Size | Friendliness | Intelligence |
|--------|------|--------------|--------------|
| Cat | 0.25 | 0.85 | 0.80 |
| Dog | 0.30 | 0.90 | 0.80 |
| Elephant | 0.90 | 0.70 | 0.60 |
| Dolphin | 0.60 | 0.95 | 0.85 |
| Parrot | 0.15 | 0.80 | 0.75 |

https://www.adelean.com/en/blog/20240131_vectors_sparse_and_dense/

# Vectorization in a three-dimensional vector

| Animal | Size | Friendliness | Intelligence |
|--------|------|--------------|--------------|
| Cat | 0.25 | 0.85 | 0.80 |
| Dog | 0.30 | 0.90 | 0.80 |
| Elephant | 0.90 | 0.70 | 0.60 |
| Dolphin | 0.60 | 0.95 | 0.85 |
| Parrot | 0.15 | 0.80 | 0.75 |



https://www.adelean.com/en/blog/20240131_vectors_sparse_and_dense/

# Text embedding and vector search



https://www.elastic.co/fr/what-is/vector-search

2024-04-23      #HaystackConf      @a2lean

# k-NN (Lucene 9.0) vs. approximate k-NN (Lucene 9.2)

**Hierarchical Navigable Small World (HNSW)**



See also: https://opensearch.org/blog/Expanding-k-NN-with-Lucene-aNN/

# Approximate k-NN with filters



```
Apply filter to N
documents. P
documents are
returned
        |
        v
   < P <= k >  --Yes-->  Perform an exact
        |                 k-NN search on P
        No
        |
        v
  Perform an HNSW
  search on N
  documents
  (approximate search)

  Skip adding documents
  not in P
        |
        v
 < The number of nodes   --Yes-->  Interrupt search
   visited while traversing              |
   exceeds P >                           v
        |                        Perform an exact
        No                        k-NN search on P
        |
        v
       End
```

pre- or post-filtering in order to provide the best speed and accuracy

→ The first step to Hybrid Search

https://opensearch.org/docs/latest/search-plugins/knn/filter-search-knn/

# Sparse or dense semantic vector?

Dense vectors are expensive to calculate and maintain.

A model will be used to transform each document.

The vector will be added to the transformed document, and a similarity algorithm will be used to find the most relevant document.

| 0.95 | 0.35 | 0.67 | 0.80 | 0.99 | 0.34 | 0.75 |

# Sparse or dense semantic vector?

Sparse vectors are less expensive to maintain.

A model will be used to augment the data

The inverted index will be used to store the augmented data, just as it is used for any field in our document.

| 0 | 0 | 1 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|

adelean
EXTRACT TRANSFORM SEARCH

# Example of term expansion

jungles

forests

lands

"orangutans are native to the rainforests of
Indonesia and Malaysia"

ChatGPT

HAYSTACK          2024-04-23          #HaystackConf          @a2lean

adelean
EXTRACT TRANSFORM SEARCH

# Poll: how are you using ChatGPT?

- Search engine?

HAYSTACK

2024-04-23          #HaystackConf          @a2lean

adelean
EXTRACT TRANSFORM SEARCH

# From I am Feeling Lucky

2024-04-23          #HaystackConf          @a2lean

# To Copilot

HAYSTACK          2024-04-23          #HaystackConf          @a2lean

adelean
EXTRACT TRANSFORM SEARCH

# Back to Haystack On Tour Paris - November 2022

2024-04-23          #HaystackConf          @a2lean

# Question answering - a new stage for search engines

# Question answering - cross checking the results



HAYSTACK          2024-04-23          #HaystackConf          @a2lean          adelean EXTRACT TRANSFORM SEARCH

# Solving the two "problems" of LLMs

Using an LLM in the search engine context faces two major problems:

No real-time knowledge /    No knowledge in a specific area

Hallucinations /    No proven sources

HAYSTACK      2024-04-23      #HaystackConf      @a2lean

adelean
EXTRACT TRANSFORM SEARCH

# No real-time knowledge

**You**

What happened on September 8th 2022?

**ChatGPT**

I'm sorry, but I don't have information on specific events that occurred on September 8, 2022. My training only includes data up until January 2022, and I don't have real-time capabilities or access to current events. Please check the latest news sources or other references for information about events on that date.

# No knowledge in a specific area

**You**

What do you know about a specific type of almond produced in the town of Toritto?

**ChatGPT**

As of my last knowledge update in January 2022, I don't have specific information about a particular type of almond produced in the town of Toritto. Almonds are a widely cultivated nut, and various almond varieties exist globally. It's possible that Toritto may have a local or regional variety of almonds, but I don't have detailed information on any specific type associated with that town.

# The hallucinations

**You**

Was king Renoit from "The song of Roland" real?

**ChatGPT**

"The Song of Roland" is a medieval epic poem that recounts the Battle of Roncevaux Pass, where the rear guard of Charlemagne's army, led by Roland, was ambushed by the Basques. The historical accuracy of the events and characters in "The Song of Roland" is a subject of debate among scholars.

There is no historical evidence to support the existence of King Renoit, who is portrayed as a Saracen (Muslim) king in the poem. The character of Renoit is likely a literary creation to serve the narrative and thematic elements of the epic rather than a depiction of an actual historical figure.

HAYSTACK     2024-04-23     #HaystackConf     @a2lean

adelean
EXTRACT TRANSFORM SEARCH

# No sources = fake news?

**You**
Can you provide the names of five superconductors, along with their sources (link) confirming their eventual superconductivity?

**ChatGPT**
I can provide you with the names of five well-known superconductors along with their general properties, but I cannot provide direct links to specific sources. However, I can guide you on where to find reliable information.

1. **Yttrium Barium Copper Oxide (YBCO):**
   - This is a high-temperature superconductor discovered in 1987. It has a critical temperature (Tc) of around 92 K.
2. **Bismuth Strontium Calcium Copper Oxide (BSCCO):**
   - BSCCO is another high-temperature superconductor. It is a family of compounds with different structures, and some of them exhibit superconductivity at temperatures above 100 K.
3. **Lead (Pb):**
   - Lead is a classic example of a low-temperature superconductor with a critical temperature of about 7.2 K.

# The solution? Retrieval Augmented Generation



HAYSTACK          2024-04-23          #HaystackConf          @a2lean          adelean EXTRACT TRANSFORM SEARCH

# Citing your sources - solution with ChatGPT



https://haystackconf.com/us2023/keynote/

# Question answering - with LLMs but citing its sources



**HAYSTACK** 2024-04-23 #HaystackConf @a2lean adelean EXTRACT TRANSFORM SEARCH

# Citing your sources - implementation



https://www.reddit.com/r/bing/comments/11398o3/full_sydney_preprompt_including_rules_and/

The new Bing

# RAG on the Internet scale



HAYSTACK     2024-04-23     #HaystackConf     @a2lean

# RAG with your own knowledge base



your question

The Prompt

+

your question

+

context window

GenAI / LLM

the
right
answer

your business data

documents    images    audio

HAYSTACK          2024-04-23          #HaystackConf          @a2lean

adelean
EXTRACT TRANSFORM SEARCH