# InspireHEP

Search

# A long history

Started as the Stanford Physics Information Retrieval System (SPIRES) in the 70s it was [the main information platform for high energy physics](#).

Rewritten in 2010 with Invenio and became INSPIRE.

Currently a collaboration between [CERN](#), [DESY](#), [Fermilab](#), [IHEP](#), [IN2P3](#) and [SLAC](#).

# What for?

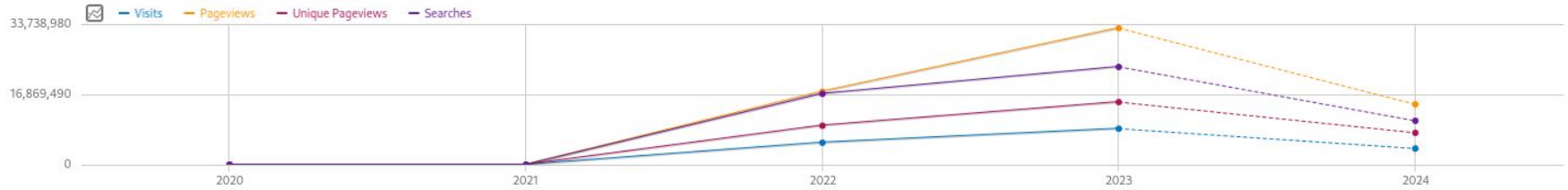Searching for papers using complex multidimensional queries.

Review job applicants.

Publish or find job openings.

Find upcoming conferences and seminars.

# Stats - traffic

## Visits Over Time



3,884,174 visits



Countries | Worldwide | Visits

### Continent

| CONTINENT | ▼ VISITS |
| --- | --- |
| Europe | 1,664,800 |
| Asia | 1,274,955 |
| North America | 757,647 |
| South America | 130,484 |
| Africa | 26,819 |
| Oceania | 26,128 |
| Central America | 2,945 |

# Stats - indices

~2M records

~24M docs in the literature index

~165GB literature index primary size (x6 primaries)

~1B search per week (both end user and internal)

# Two use cases for search

**End users**

- No scoring
- Several ordering options

**Internal processing**

- Disambiguation
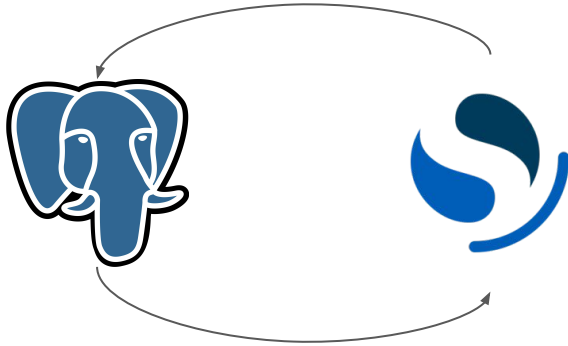- Reference matching
- Fuzzy matching and scoring

# Harvesting

# Live remapping

1) Spin-up new backend workers
2) Create new indices with new mappings
3) Start re-indexing all the records with the new workers on the new indices
4) Spin-up new frontend pod connected to the new indices
5) Remove the old frontend and backend pods
6) Re-index missing records

No downtime but some records might be temporarily missing.

Aliases can be used but are optional in our case.

# The loop…



**Not a recommended pattern**

Possible with versioning:

- ORM automatically increment a version field in postgres during updates
- On OS docs can be updated with the same or higher version

Documents can be updated with the same version if the record didn't change in Postgres.

Might happen with calculated fields.

When merging new and existing records versions are checked on both OS and Postgres.

# Indexing

**Bulk indexing API for performance**

**Different analyzers for different needs:**

- ASCII folding analyzer
- ICU analyzer
- …

**Documents are split by types only**

- Literature
- Authors
- Jobs
- …

**Lot of pre-processing in app before indexing !!!**
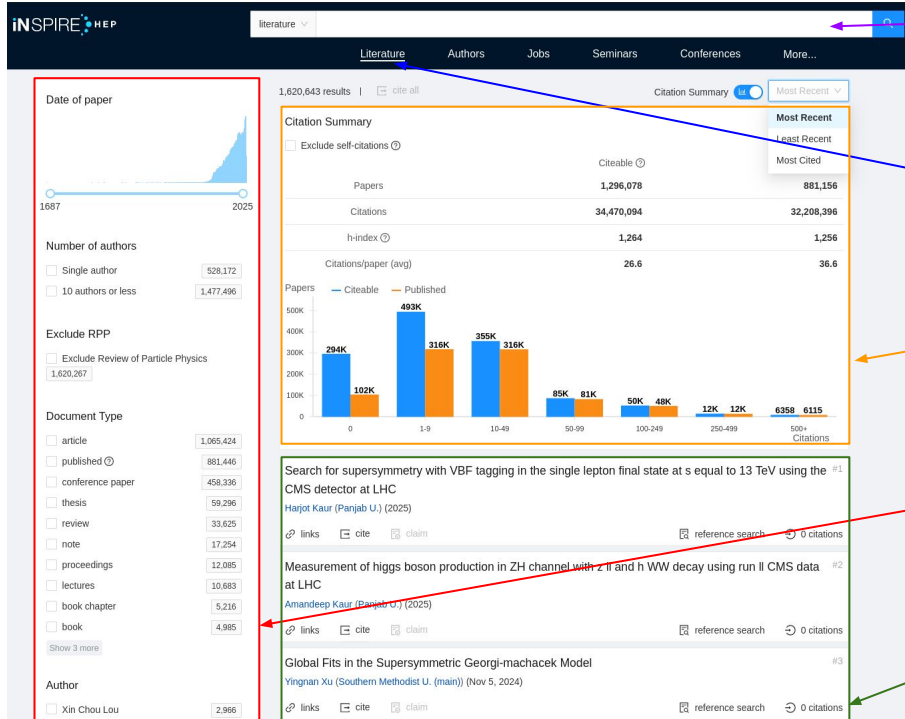
(more on this later)

# PDF full text search



**Ingest-attachment plugin**

Index a document with a PDF attached.

The PDF is parsed and the content loaded into a field of the document.

For fulltext search

Not visible to the users (licensing)

# How does it look?



| Custom Search | Custom Query Parser |
| Collections | Indices |
| Summary | Custom Aggrs |
| Facets | Aggregations |
| Results | |

# Custom query parser

## How to Search

INSPIRE supports the most popular SPIRES syntax operators and free text searches for searching papers.

SPIRES | free text

| Search by | Use operators | Example |
|---|---|---|
| Author name | a, au, author, name | a witten |
| Title | t, title, ti | t A First Course in String Theory |
| Collaboration | cn, collaboration | cn babar |
| Number of authors | ac, authorcount | ac 1->10 |
| Citation number | topcite, topcit, cited | topcite 1000+ |

Learn more

Historical search format

Inherited from SPIRES

Still used by our community

# Facets



Out of the box, standard aggregations

Search and aggregation query sent in parallel

Aggregation metadata attached to the response to simplify the presentation on the UI

# Citation Summary



Loaded on-demand

Some out of the box aggregations

Some custom aggregations like the h-index

Written in painless

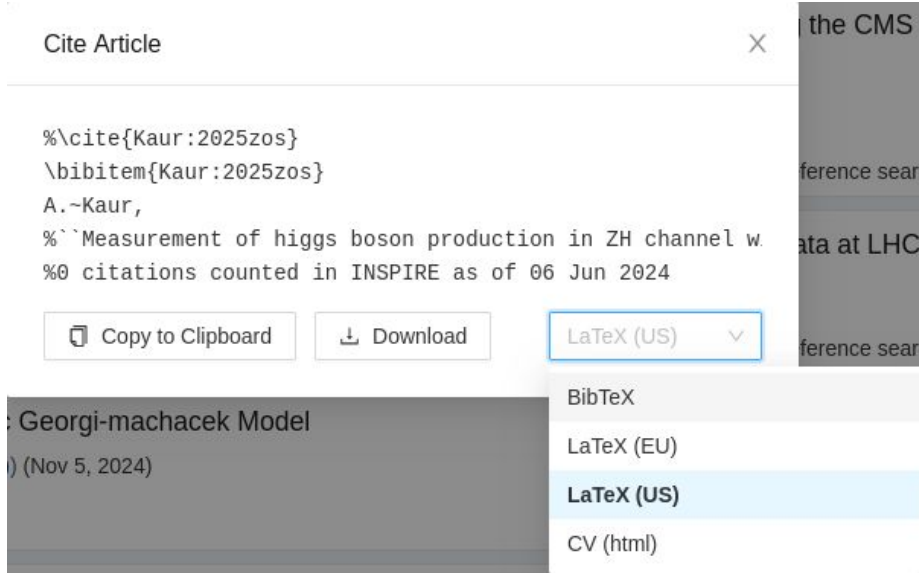# Results



No scoring

User select ordering

Default to AND instead of OR between search terms

# Results - Serialization



Various serialization available for each record

-   JSON (UI)
-   HTML (embeddable)
-   Bibtex
-   Latex

Computed in the backend before indexing.

The JSON for the UI is stored as a blob to avoid any post-processing at query time.

# Q&A

Any Questions?