# A new textual search engine to discover VizieR catalogs

**CDS** CENTRE DE DONNÉES ASTRONOMIQUES DE STRASBOURG

Textual search is a part of the VizieR indexation which completes the position indexation and the keyword search resulting from SQL queries. This new capability extends the capacity of the current engine with a natural language approach like in Google, or bumblebee (ADS). The new version – still in alpha – uses the Elasticsearch engine, an Open Source Software that indexes documents with a grammar and a textual search analyser.

The query supported is a NO-SQL language including strict or fuzzy search and available through an HTTP RESTful API. Then, a fine configuration adapted to the different data (authors, DOI, abstract, date...) is needed to improve the indexation.

Resources indexed are the ReadMes which describe the VizieR catalogs. A ReadMe is a structured ASCII file containing the basis metadata : authors, title, keywords, abstracts but also tables and column descriptions. We explain the new implementation from the data origin to the final users.

## VizieR will use Elasticsearch for textual search
in addition to the other searches like search by position

**elasticsearch**

### A word on Elasticsearch operation
Elasticsearch, based on Apache Lucene engine is allowing to implement a flexible search engine easily. In this purpose there are three phases to follow : index configuration, data ingestion and finally search.
The first part is to configure an index by choosing the way data should look like and how they will be treated. Once done, you need to format your data into JSON, the data format for Elasticsearch communication. Then you are ready to ingest them in Elasticsearch engine. Finally you can use the Elasticsearch API to query the index.

## Fig.1: A new interface
### Search for Year:>2010 AND exoplanet AND bibcode:*MNRAS*



### A new alpha interface to access the « VizieR bazar »

The fig.1 on the right shows an example of Elasticsearch usage. It gives you access to a complete grammar for requests, allowing to use operators, fuzziness (levenshtein distance), phrase analysis, proximity (number of words between words of a phrase) and other. All that you can combine the way you want.

The chosen vocabulary is based on ADS field naming : author, title...
Now it is also possible to search on doi, orcid, year of publication.
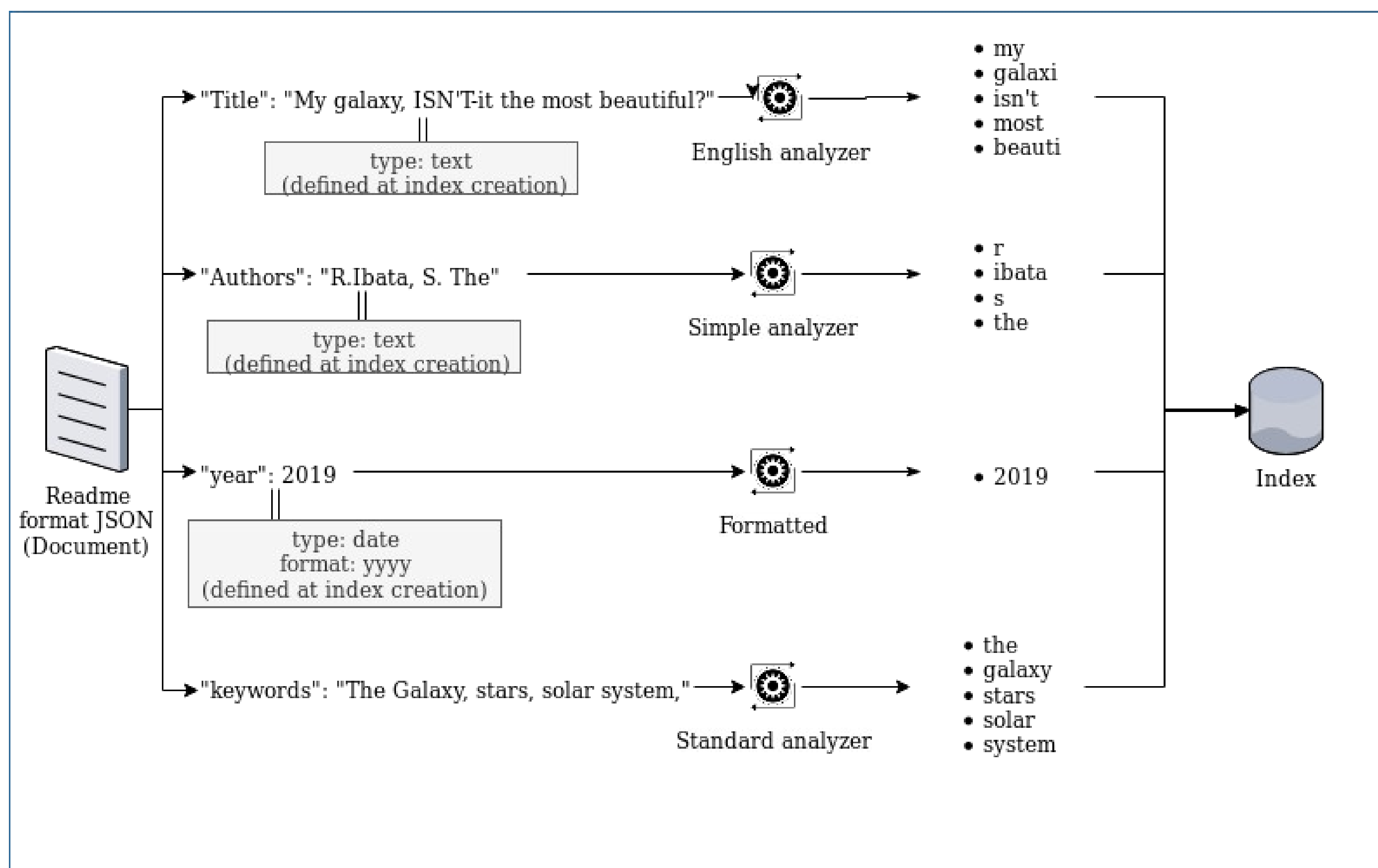
The alpha version is available here :
http://cdsarc.u-strasbg.fr/viz-bin/cat2

### Analysis
On fig.2 you can look at an example of a document analysis, here a Readme. It shows 4 fields treated differently for the indexation and the effective outcome that will be indexed.

### Current work
Elasticsearch ingestion and search are ready for VizieR in local CDS architecture and will be proposed to the outside (VizieR mirrors and other workflows).
Currently, the mirroring architecture has been tested locally in a docker environment.

### Perspectives
In the future, Elasticsearch for VizieR may be used to populate the VizieR catalogs in the VO registry. Some tests are planned to use Elasticsearch to feed the chatbot experimented to access the CDS servers.

fig.2: ElasticSearch documents analysis (indexation)

**Bisch.Y, Landais.G, Schaaff.A**
yann.bisch@astro.unistra.fr

**VizieR**

**CNRs** UNIVERSITÉ DE STRASBOURG