# NLP-chatbot R&D at CDS. VO standards, AI impact, etc.

André Schaaff, Thomas Boch, Sébastien Derriere

Centre de Données astronomiques de Strasbourg
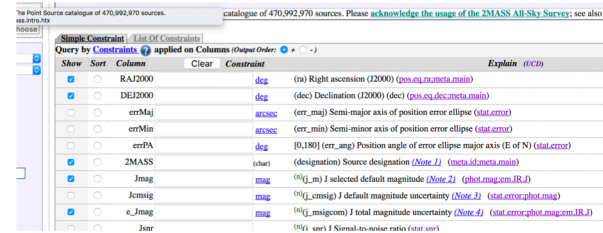
IVOA Malta, Valetta Campus, 14-17 November 2024

# ☐ Purpose

- Feedback and status of a long term R&D action at CDS

- A study around Natural Language Processing applied to the astronomical data services querying

- It was presented step by step in various IVOA groups

- Taking now into account the AI vawe…

# Why NLP ?

- Since Web advent we are used to deal with webpages containing forms (parameter fields, checkboxes, etc.)



- Nowadays we use the voice (gesture should join for multimodal requests) to interact with several devices with voice answer or a displaying

- An interesting field of investigation with both technical and human concerns
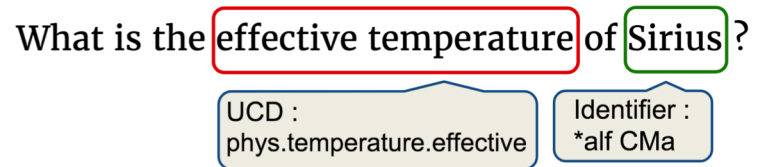
# Webpages / Forms pro cons

- Forms are not (yet) old fashion and have still advantages
  - Easy and quick for experimented users
  - easy to translate to database queries
  - Result rendering "easy" to design

- On the negative side :
  - not intuitive for new or occasional users
  - Thus we use often only a (small) part of the functionalities

- Understanding / Knowledge / Learning efforts on the user side

# NLP pro cons

- A friendly (Inclusive if disabilities) way of interaction
  - More digestible than forms
  - Just a prompt or a voice recognition
  - Result displaying can be similar than for forms or a speech synthesis

- On the negative side :
  - possible frustration for expert users if the result is too imprecise
  - too long chat if the question was vague or ambiguous

- User request Understanding and translation to query the services is on the service side

# Pragmatic and best effort approach

- Learning about NLP (basis, tools, examples, …)
- Starting by building a set of typical queries

- We had already a lot of helpful data and tools
  - authors in Simbad, VizieR
  - missions and wavelengths in VizieR
  - DJIN to recognize identifiers in a text
  - UCDs
  - a name resolver



What is the effective temperature of Sirius ?

UCD : phys.temperature.effective

Identifier : *alf CMa

# And a useful friend : the Virtual Observatory

- We used VO standards like TAP/ADQL (DAL WG), UCDs (Semantics WG), … ,and tools based on ML (and now AI), a link with KDD IG.

- The VO was a master piece which helped us to translate NL into queries to our own services. By the way all services implementing VO standards could be included.

List the QSOs at Z > 6.

Simbad, TAP query

```
SELECT main_id, oid, rvz_redshift
FROM basic
WHERE otype = -14680064 AND rvz_redshift > 6;
```

What is the effective temperature of Sirius ?

VizieR, TAP Query

```
SELECT "VI/137/gum_mw".Teff
FROM "VI/137/gum_mw"
WHERE 1 = CONTAINS(POINT('ICRS', "VI/137/gum_mw"."RAJ2000",
"VI/137gum_mw"."DEJ2000"), CIRCLE('ICRS', 101.287155333,
-16.716115861, 20/3600.)) ;
```

# On the  way to a chatbot

- A first work with the Stanford NLP library was far from a Chatbot  as it must be able to chat with the user otherwise it is just like  a command line

- Google Dialogflow brought us this chat capability
  - Based on Machine Learning
  - Training, tuning, improving features

- Query understanding  was then delegated to Dialogflow

# Switching from coding to training to understand user intents
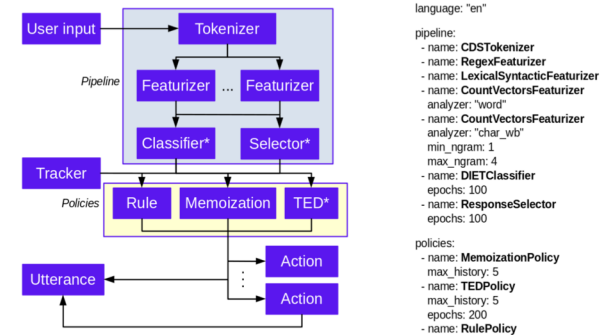
# Examples of sentences linked to get_measure

# Remark: External tools

- With Dialogflow we delegated a larger part of the tool than with the first approach but not the link to the services.
- The link to the services and de facto to the data is the key part and our development was always  designed to switch to other NLP tools.
- We switched to an Open Source framework, Rasa

# How does Rasa work ?

- … Rasa prediction model training is based on Tensorflow

- Compared to Dialogflow the configuration is not really user-friendly

  – Items of the architecture have to be written
  – For example a CDSTokeniser

# Current version based on Rasa framework



VOA Malta, KDD - GWS Joint Session, 16/11/2024

# The Day after

- Work with Rasa was continued (cleaning, refactoring, etc.) but with the advent of chatGPT we decided to make a break and to spend time on the study of its impact, a help for the existing chatbot ? or a way to reshape it completely ?

- We had a look at several AI tools and we decided to test with the OpenAI API which cost model (credits to buy and not a monthly subscription ) is possible for us

# First Experiments with OpenAI

- OpenAI has a functionnality called  assistant which allows the definition of assistants inside an application, able to answser a request using models, tools ans data specified by the developper

- A first prototype was defined to learn  Simbad table structure and ADQL to become able to propose or verify  TAP queries to Simbad.

# Examples

**User**

give me the TAP request of the position of Vega

**TAP Query**

To obtain the TAP (Table Access Protocol) request for the position of Vega, I will use the "basic" table, which contains general data about astronomical objects, including their coordinates (right ascension and declination) 【7†source】 . The TAP query will involve selecting the right ascension (RA) and declination (Dec) for the object identified as Vega. Let me construct the TAP query for you.

**TAP Query**

The TAP query to retrieve the position of Vega from the Simbad database would be:

```
1  SELECT ra, dec
2  FROM basic
3  WHERE main_id = 'Vega'
```

This query selects the right ascension (ra) and declination (dec) from the "basic" table, where the main i (main_id) is equal to 'Vega'.

**User**

this not work because vega is not the main id. You must join with the indentifiers table

**TAP Query**

You are absolutely right. To obtain the position of Vega, I will need to join the "basic" table with the "ident" table using the object internal identifier (oidref).

The TAP query to retrieve the position of Vega from the Simbad database with the necessary join would be:

```
1  SELECT b.ra, b.dec
2  FROM basic AS b
3  JOIN ident AS i ON b.oid = i.oidref
4  WHERE i.id = 'Vega'
```

This query selects the right ascension (ra) and declination (dec) from the "basic" table and joins it with the "ident" table on the object internal identifier, then filters the result to only include the entry with the identifier 'Vega' in the "ident" table.

# Ongoing experiments

- We have a mine of information not directly used in the queries to the services like the VizieR Readme files



LLMs are efficient for textual content
The aim is to train with this additionnal data to detect the astronomical object names in the Readme files

Files of valid, invalid names to train, etc.

# ☐ remark

- Testing / working with AI frameworks is often nearer to cooking than to programming
- best recipe to find
- Trying with different models
- Training with a good set of data (also a set with errors)
- Improving with other sources or various data formats for a better understanding by the models
- Etc.

# Conclusion

- We are testing and improving step by step to interact with data services differently, keeping in mind the gap between a prototype and a production service

- .. And adapting us to new tools/frameworks/paradigms; so as usual in IT

Providing stable and sustainable new tools/services benefiting from AI is an interesting challenge

Questions ?