# Automatic recognition of object names in literature

**C. Bonnin**[1], S. Lesteven[1], A. Oberto[1], S. Derriere[1], F. Genova[1]

[1] CDS – Strasbourg

## Bibliographic code :

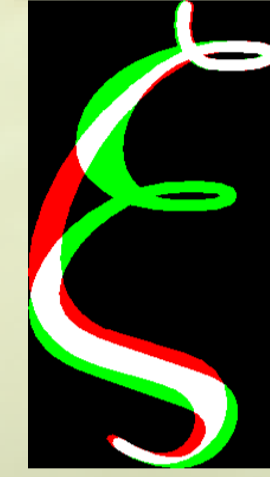**2007MNRAS.374..176L**

Connection to the Journal Bibliographical Service to retrieve the PDF document corresponding to the bibliographic code.

## Raw PDF document

• Text extraction using PDFBox
• Pictures of special characters using JPedal
• Recognition of graphic symbols
   (here comparison between greek letters
• $\zeta$ and $\xi$ )

## Extracted text

•

Dictionary of nomenclature    SIMBAD
•
•
•

• Searching for object names using :
• the formats from the Dictionary of Nomenclature of Celestial Objects
• usual names such as Aldebaran or The Crab given by SIMBAD
• the variable star names based on constellation names

Use of the Weka Machine Learning engine to help detecting the false positives among the object names.

WEKA The University of Waikato

The object names are highlighted in the original document and displayed by Acrobat Reader.
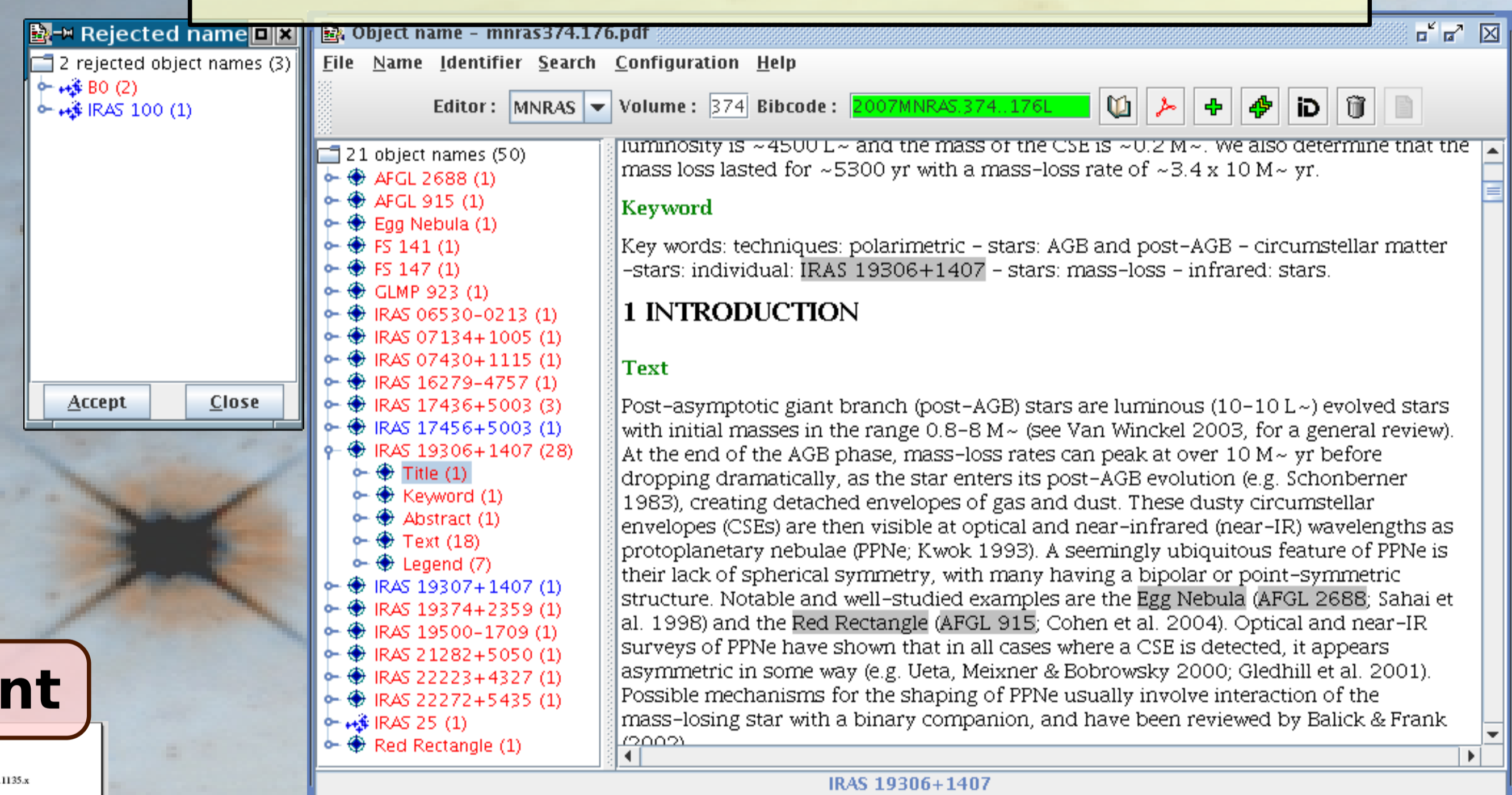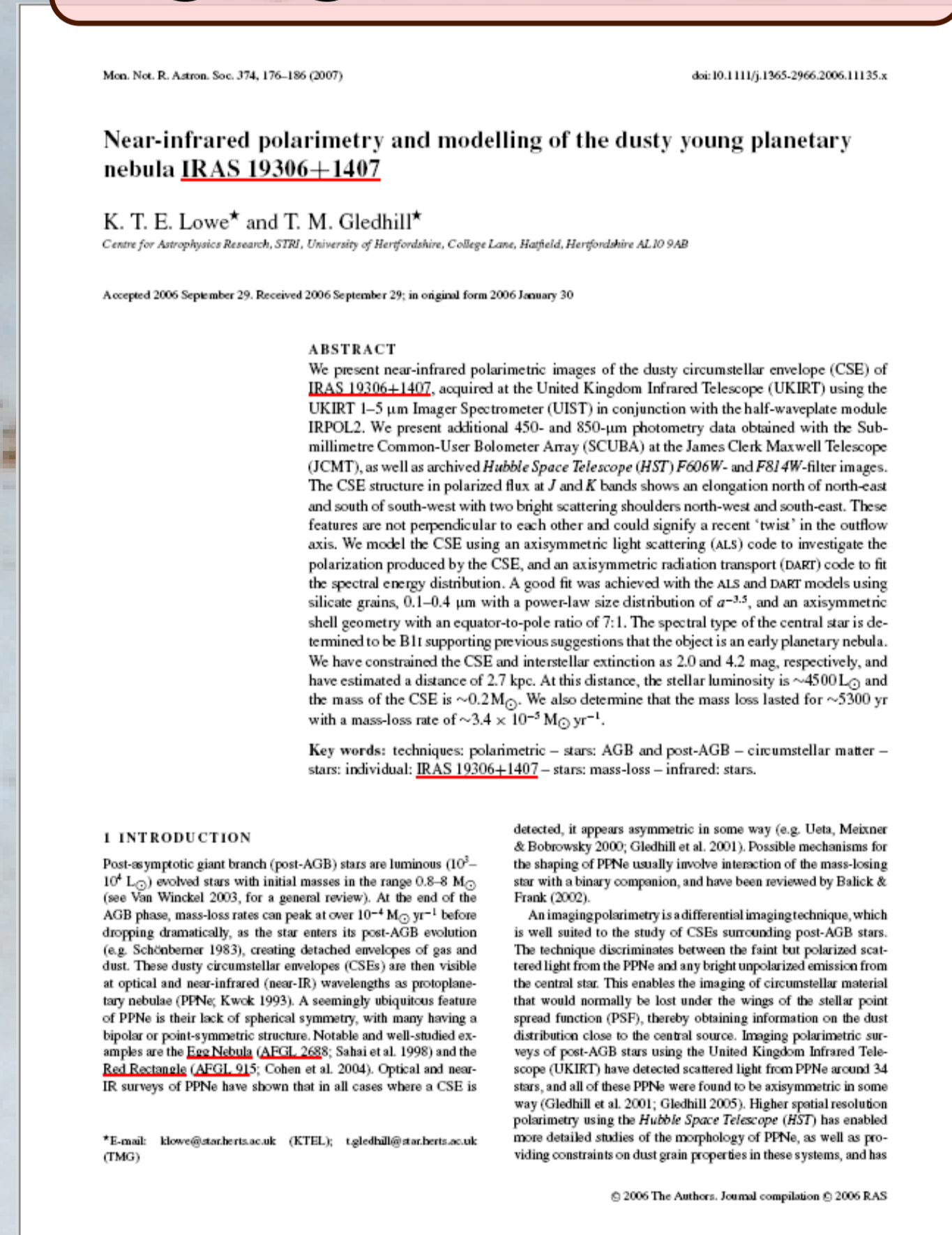
## Highlighted document



## Today :

SIMBAD is a database of astronomical objects that provides (among other things) their bibliographic references in a large number of journals. Currently, these references have to be entered manually by librarians who read each paper.

## • Objectives :

Detect object names in published articles and propose them for validation by the experts.
This includes :
• Search for possible object names in PDF documents
• Highlight the names in the documents without changing the presentation
• Retrieve information about the objects in SIMBAD
• Display the possibilities and let the librarian make the final choice
• Gather some additional information such as the number of occurrences and the positions in the document
• Enter the references and additional information in SIMBAD

A graphic user interface written in Java Swing can perform all these operations.



Each object name is checked for existence in SIMBAD and the object description is retrieved. A librarian makes the final validation.

## Identifiers list

| Name | Position | Nb |
|---|---|---|
| AFGL 2688 | text | 2 |
| AFGL  915 | text | 2 |
| FS 141 | text | 1 |
| FS 147 | text | 1 |
| IRAS 06530-0213 | text | 1 |
| IRAS 07134+1005 | text | 1 |
| IRAS 07430+1115 | text | 1 |
| IRAS 16279-4757 | text | 1 |
| IRAS 17436+5003 | text | 4 |
| IRAS 19306+1407 | text, title, keyword, abstract, legend | 30 |
| IRAS 19374+2359 | text | 1 |
| IRAS 19500-1709 | text | 1 |
| IRAS 21282+5050 | text | 1 |
| IRAS 22223+4327 | text | 1 |

• The information is entered directly into SIMBAD :
• Object identifier
• Name as it is written by the author
• Position of the object citation in the text (title, abstract, keyword, table, ...)
• Number of occurences

## Update

SIMBAD

## • Results :

• Improvement and automatic verification of the Dictionary of Nomenclature of Celestial Objects
• More verification of data entered in SIMBAD
• Documents treated more quickly and more exhaustively
• Help for the librarians to concentrate on the added value of their work
• Beta version currently in use

• More information will be made available to find the most relevant papers in the object reference lists (number and position of occurrences).