# CDS homogenisation of metadata from publishers
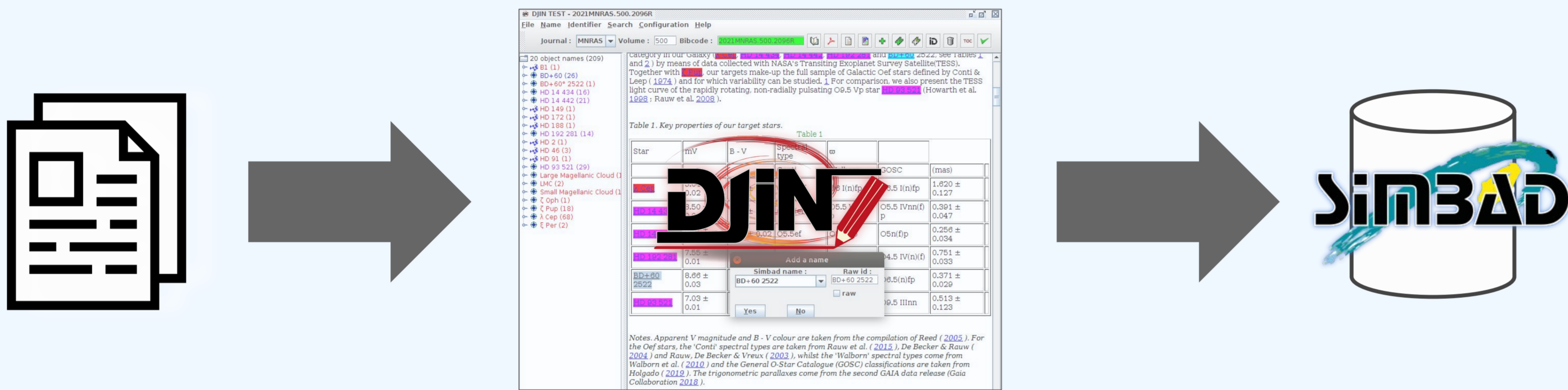
**Grégory Mantelet**
gregory.mantelet@astro.unistra.fr
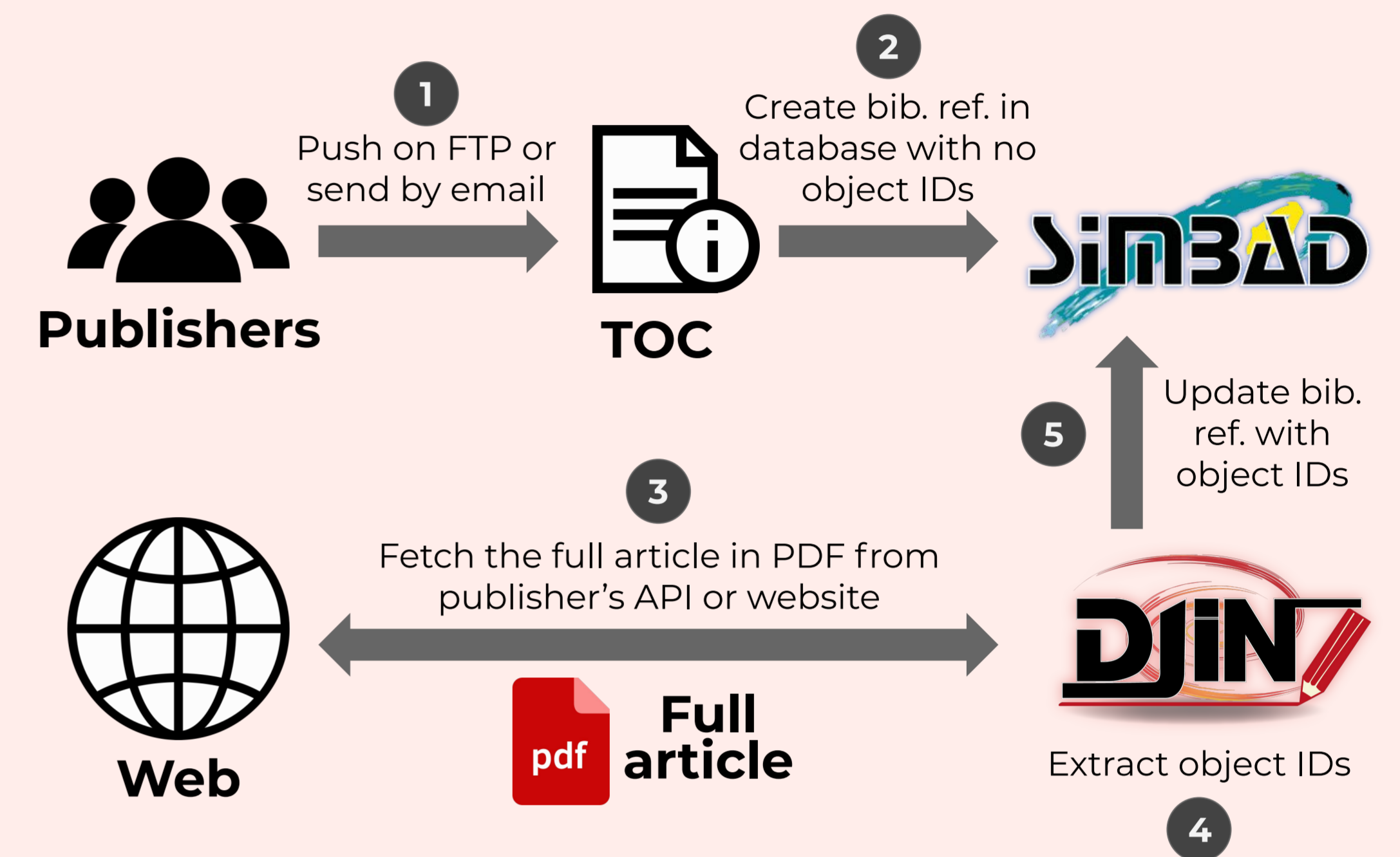**Oberto A., Neuville M., Lesteven S., Allen M.**

## What purpose?

- **Articles analysed**...
- ...for astronomical **object names extraction**
- ...by **DJIN** (a semi-automatic software)
- ...to update the objects stored in the **SIMBAD database**
- ...with new identifiers, positions
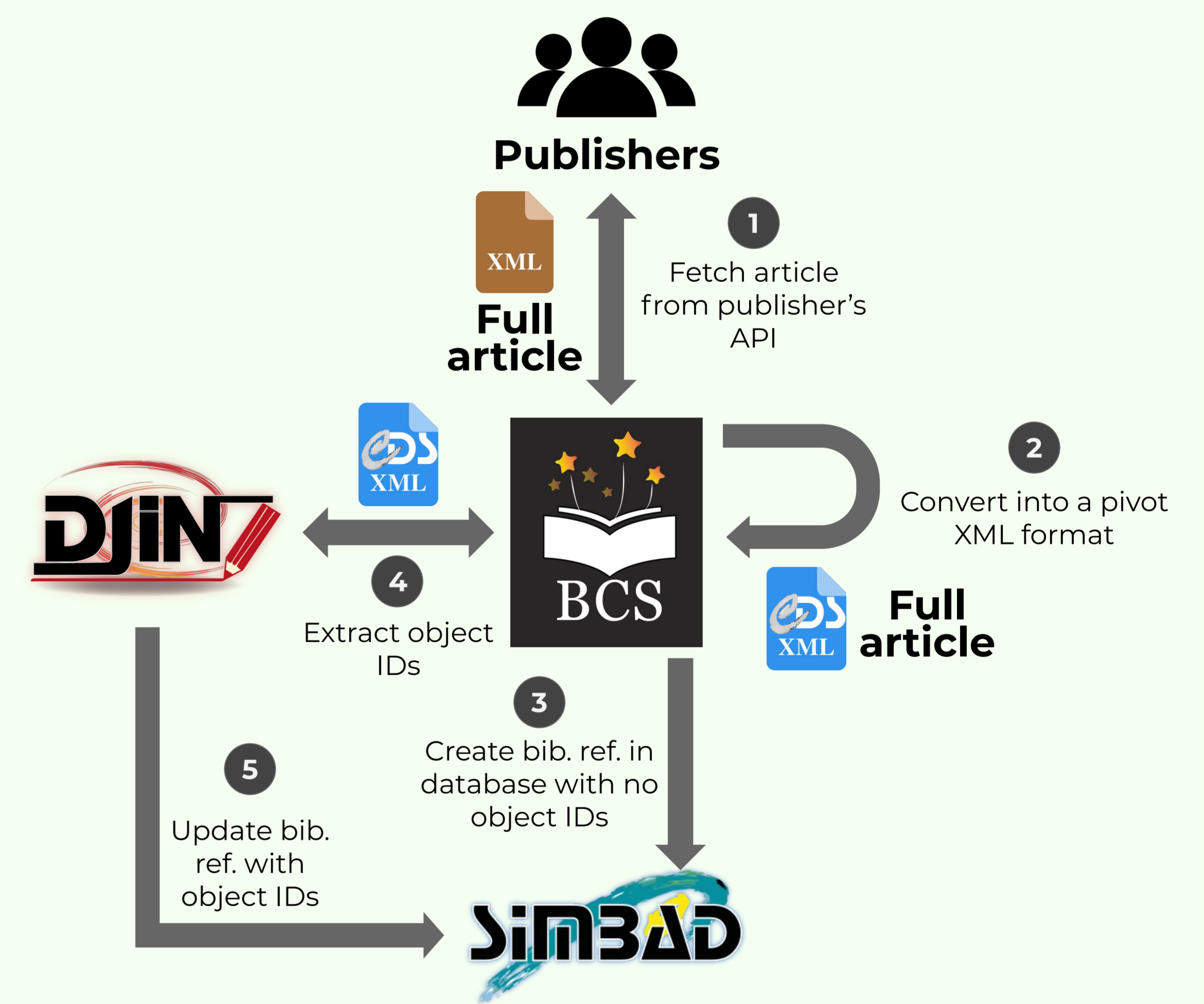- ...and **links between objects and their bibliographic references**.



## Now



## What change?

| | **Now** | **In progress** |
|---|---|---|
| **Format** | **PDF** <br> • not reliable for text extraction but for display and print, <br> • different from one publisher to another, <br> • very sensitive to change | **XML** <br> • easy text extraction (XPath) <br> • easy conversion into other formats (XSLT) <br> • if following a **standard**, XPath+XSLT reusable (JATS) |
| **Files** | **2 files** <br> with a different way to get each of them depending on the publisher | **Only 1 file** <br> with a different way to get it depending on the publisher |
| **Workflow** | **Completely manual + Scripts** (Bash/Python) <br> get and process the TOC + get the PDF article from the Web | **Semi-automatic + Internal Web Site** <br> fetch+convert on demand with the BCS when available |
| **Get** | **TOC: Push** <br> (pushed by publishers on our FTP or by email) <br> **PDF: Push** on our FTP <br> or **PDF: Pull** from web site | **Pull** <br> (get XML from publishers) <br> *As much as possible, but not yet possible with most publishers* |
| **Reusable?** | **No** | **Yes** <br> *e.g. for table extraction (VizieR)* |

## In progress



### BCS
**Bibliographic Center Supervisor**
**(Internal CDS Web Site)**



## But...

- **Not all** journals provide their articles **in XML**
- If existing, the **XML** is **not always accessible for us**
- If existing, **often not standard** (though there are standards like **JATS**)
- **Missing** Pull **Service or API** to fetch XML articles
- **Missing notifications** about **complete** new volumes/issues