

The diagram illustrates a six-step process flow, divided into two parts:

- SKOL Part II (Green):**
 - 01 Parse:** The first step, supported by **PySpark Automation**.
 - 02 Classify:** The second step, also supported by **PySpark Automation**.
 - 03 Remove Extraneous content:** The third step, supported by **PySpark Automation**.
- SKOL Part I (Grey):**
 - 04 Generate Taxon:** The fourth step.
 - 05 * Embed Taxon Content In Search:** The fifth step.
 - 06 * Provide Semantic Search Interface:** The final step.

A bracket labeled **SKOL Part I** spans steps 04, 05, and 06. A bracket labeled **SKOL Part II** spans steps 01, 02, and 03. The **PySpark Automation** bar is positioned below steps 01, 02, and 03.

Annotated 244 issues from old issues of 3 freely available journals:

- Persoonia, Mycotaxon, and Mycologia
- from mykoweb.com

[@Glomus cubense Y. Rodr. & Dalpé, sp. nov.*Nomenclature*]
[@Plates 1-2#Figure*]
[@Mycobank MB561650#MB*]
[[@Sporocarpia ignota. Sporae singulares vel aggregatae. Sporae hyalinae vel luteolae, ovoideae, vel irregulares, 20-48 x (24-)54-82 µm, raro globosae, 24-54 µm diam. Sporae tunica stratis duobus; stratum exterius hyalinum, subflexuosum, ad 0.8-1.0 µm crassum, stratum interius hyalinum vel luteolum, 0.8-1.7 µm crassum, rigidum. Hyphae sustentines hyalinae, rectae, subcylindricae vel subinfundibuliformes, 4-10 µm crassae ad basin sporae. Porus apertus, raro septo curvo clausus cum strato interiore. Mycorrhizas vesicular-arbusculares formantes.#Description*]

[illegible]

- Based on Dr. Draft's SOTA Literature Search from the Auton Lab.
- Uses a Large Language Model (SBERT) to create a space of all nomenclature/description sets.
- Provide a description on the command line which is embedded into the space.
- Returns the closest matching descriptions by cosine similarity.
- Should be publicly available at <https://synoptickeyof.life> by end of June 2025.

- over 1000 unannotated documents to ingest
- Parses OCR output into paragraphs
- Classifies paragraphs into Nomenclature, Description, Misc-exposition.
- Distributed pipeline to allow parallel processing.

- Uses an LLM (gemma) to extract features and their values from each description
- Annotate the resulting json with TF-IDF scores for every feature
- Basis for building a synoptic key User Interface

- Web bot to collect open content articles
- Automated OCR pipeline
- Provenance annotation
- Front end

Collaborators:
Balasi
Caspers
Murphy
Osuga
H.P. Yarroll

coming soon: <https://synoptickeyof.life>
<https://github.com/piggyatbaqaqi/skol>
<https://github.com/piggyatbaqaqi/dr-drafts-mycosearch>
 IST664 Synoptic Key of Life
 IST718 Synoptic Key of Life II
<https://github.com/autonlab/dr-drafts-sota-literature-search>

This poster:

