PREDICTS 2 ETL Process

1 CONTENTS

Sum	mary	1
Usef	ul links	1
Extra	act names	2
3.1	Load unresolved names from database	2
3.2	Load manually resolved names	2
3.3	Merge streams	2
Clea	n + transform	2
l.1	String cleaning	2
1.2	Build GBIF query	3
Quei	ry GBIF	3
5.1	GBIF species match query	3
5.2	Transforming results	3
Evalu	uate GBIF matches	4
5.1	Handling unmatched records	4
5.2	Handling matched records	4
5.3	Check match confidence	5
Writ	e out results	5
7.1	Resolved names	5
7.2	Unresolved names	7
	Usef Extra 3.1 3.2 3.3 Clea 4.1 4.2 Que 5.1 5.2 Evalu 5.2	3.2 Load manually resolved names 3.3 Merge streams Clean + transform 3.1 String cleaning 3.2 Build GBIF query 3.1 GBIF species match query 3.2 Transforming results 5.2 Transforming results 5.1 Handling unmatched records 5.2 Handling matched records 5.3 Check match confidence Write out results 7.1 Resolved names

1. SUMMARY

This document records the ETL process used by NHM Informatics to resolve taxonomic names gathered from biodiversity papers as part of the PREDICTS2 project. The GBIF species endpoint is used to resolve straightforward names (including synonyms) and to identify names in need of manual resolution by an appropriately qualified human. The steps laid out in the rest of this doc should be read in conjunction with ETL diagram.jpg.

2. USEFUL LINKS

ETL platform: Pentaho Data Integration (Community edition v8.3)

GBIF Species API (as of 01/04/2019 – 30/09/2020): https://www.gbif.org/developer/species

3. EXTRACT NAMES

3.1 LOAD UNRESOLVED NAMES FROM DATABASE

Unique (original species name, source name) pairs without a resolved name are extracted from dbo.predicts_2.species_record.

```
SELECT DISTINCT so.source_name, sr.species_name AS 'original_species_name',
sr.high_level_taxa AS 'high_level_taxa', sr.species_name AS 'original_species_name'
FROM species_record sr, diversity_report dr, sample_event se, site s,
study st, source so
WHERE sr.name_resolution_id IS NULL
AND sr.diversity_report_id = dr.diversity_report_id
AND dr.sample_event_id = se.sample_event_id
AND se.site_id = s.site_id
AND se.site_id = s.site_id
AND st.source_id = so.source_id
ORDER BY sr.species_name ASC
;
```

Query 3-1 - unresolved name extraction

3.2 LOAD MANUALLY RESOLVED NAMES

A locally-synced copy of needs_manual_name_resolution.xlsx is read in. It contains unique (original species name, source name) value pairs and a manually_corrected_species_name field where the PREDICTS2 team should enter any manually corrected taxonomic names we've been unable to match against GBIF.

If manually corrected species name is NULL, the record is filtered out shortly after load.

3.3 Merge streams

As each record from 3.1. is processed, its (source name, original species name) fields are checked against the records from 3.2. On match, manually_corrected_species_name is returned and attached to the record.

field	description
source_name	Name of source/study the name was found in
original_species_name	Verbatim, unresolved species name
manually_corrected_species_name	Corrected form of name

Table 1 - fields present in stream at end of extract phrase

4. CLEAN + TRANSFORM

4.1 STRING CLEANING

name_string_cleaned field is created. It contains a copy of original_species_name where ". | _ | ," characters have been replaced by whitespace. Non-standard delimiters can interfere with GBIF's fuzzy-match search. E.g.,

- Boykinia.major: match confidence 94 against (correct) genus-level result
- Boykinia major: match confidence 98 against (correct) species-level result

Occurrences of the strings ["species", "spp", "sp"] are also removed in name_string_cleaned. Use of these strings to e.g., indicate an uncertain species name impacts on the fuzzy-match search:

- Antennaria: tied match confidence 99 against genus-level result
- Antennaria spp: match confidence 74 against genus-level result

4.2 BUILD GBIF QUERY

If manually_corrected_species_name is NULL, new field name_to_use is set to the value
of name_string_cleaned. Else, name_to_use =
manually_corrected_species_name.

searchTerm is created. It contains a URL-encoded copy of name_to_use, which is appended to the end of the GBIF fuzzy match URL: <a href="http://api.gbif.org/v1/species/match?verbose=true&name="http:

field	description
source_name	Name of source/study the name was found in
original_species_name	Verbatim, unresolved species name
manually_corrected_species_name	Corrected form of name
name_string_cleaned	Cleaned version of species name
name_to_use	Version of the name that will be used to query the GBIF API
searchTerm	GBIF URL

Table 2 - fields in stream at end of clean/transform

5. QUERY GBIF

5.1 GBIF SPECIES MATCH QUERY

Multi-threaded calls to the GBIF API are made using the searchTerm URL created in the previous step (round robin distribution). Results are received in JSON and stored in a new field: result.

5.2 Transforming results

JSON record stored in the result field is flattened. At this stage, only values relating to the quality and nature of the GBIF match are extracted and appended to the record (the result field is unchanged). The current system timestamp is also added to the record

New field name	JSON path	Туре	Description
usageKey	\$.usageKey	Integer	GBIF ID of usage
matchType	\$.matchType	String	Match type (controlled vocabulary)
status	\$.status	String	Taxonomic status (controlled vocabulary)
acceptedUsageKey	\$.acceptedUsageKey	Integer	GBIF ID of accepted usage if status =
			SYNONYM
confidence	\$.confidence	Integer	Measure of confidence in top match (0-
			100)
note	\$.note	String	Field containing more detailed match
			information.

alt	\$.alternatives	String	List of alternative matches in decreasing
			confidence order. JSON.

Table 3- fields added after query stage

6. EVALUATE GBIF MATCHES

6.1 HANDLING UNMATCHED RECORDS

The value of matchType is used to identify records without a clear 'best' match against the GBIF backbone. These fall into one of three categories:

- a. no matches found: alt is NULL, matchType = 'NONE'
- b. > 1 equally good match found: note starts with 'Multiple equal matches', matchType = 'None'
- c. no 'good' matches found: alt is not NULL, matchType = 'NONE'

Records of type a. or b. are written to needs_manual_name_resolution.xlsx_when the pipe terminates. See section 7.2 for detail.

To handle records of type c. (i.e., matches exist, but the confidence level is below the GBIF quality-control confidence cut-off), the value of alt is flattened as described in section 5.2. This may comprise more than one record.

The alt result with the highest confidence is identified, extracted and fed back into the matched record process (see section 6.2). Most of these matches have low confidence and will require manual resolution, but this path was worth handling separately to allow finer control over the cutoff value.

6.2 HANDLING MATCHED RECORDS

Where a single, good match is found in the GBIF backbone, the matchtype field is examined and used to filter the results into two streams: synonyms and accepted names.

- **Synonyms:** acceptedUsageKey is used to construct a second GBIF query to retrieve the accepted name. The record is then parsed in the same step as non-synonyms, see next.
- Non-synonyms: These are flattened to extract taxon information, hierarchy and GBIF indicators:

Name	Path	Туре
scientificName	\$.scientificName	String
canonicalName	\$.canonicalName	String
rank	\$.rank	String
t_status	\$.taxonomicStatus	String
kingdom	\$.kingdom	String
kingdomKey	\$.kingdomKey	Integer
phylum	\$.phylum	String
phylumKey	\$.phylumKey	Integer
order	\$.order	String
orderKey	\$.orderKey	Integer

family	\$.family	String
familyKey	\$.familyKey	Integer
genus	\$.genus	String
genusKey	\$.genusKey	Integer
species	\$.species	String
speciesKey	\$.speciesKey	Integer
class	\$.class	String

Table 4 – taxonomic fields extracted

NB: descriptions for most of these fields are given in the GBIF API docs: https://gbif.github.io/gbif-api/apidocs/org/gbif/api/vocabulary/package-summary.html, and https://www.gbif.org/developer/species

6.3 CHECK MATCH CONFIDENCE

Records are evaluated according to their confidence measure: only records with a value >= 95 are written to the database.

The confidence measure and cut-off point were selected following testing and analysis of available options. Some 600 unmatched names were fed into the pipe while variables such as cut-off level and similarity measure (including alternative fuzzy-matching algorithms e.g., Jaro, Levenshtein, Wunsch) were modified between runs.

The F1 score of confidence \geq = 95 = 0.95, with precision of 1.

7. WRITE OUT RESULTS

7.1 RESOLVED NAMES

Accepted matches are de-duplicated and written to two tables in dbo.predicts_2_live:

a. tbl.gbif_species_reference holds taxonomic data retrieved from GBIF: one record per species/genus record. If the GBIF usage_key is already present in the table, the record is updated.

predicts_2. gbif_species_reference	Stream field	Update on duplicate key
usage_key (PK)	usage_key	Υ
scientificName	scientificName	Υ
canonicalName	canonicalName	Υ
rank	rank	Υ
taxon_status	taxon_status	Υ
kingdom	kingdom	Υ
phylum	phylum	Υ
order	order	Υ
family	family	Υ
genus	genus	Υ
species	species	Υ

class	class	Υ
kingdomKey	kingdomKey	Υ
phylumKey	phylumKey	Υ
orderKey	orderKey	Υ
familyKey	familyKey	Υ
genusKey	genusKey	Υ
speciesKey	speciesKey	Υ
updated	name_parser_timestamp	Υ

Table 5 - contents of dbo.predicts_2.gbif_species_reference

b. tbl.name_resolution holds the audit trail/results of the name resolution process:

predicts_2.name_resolution	Stream field
name_resolution_id (PK)	[auto-increment on insert]
gbif_usage_key (FK)	usage_key
original_species_name	original_species_name
manually_corrected_species_name	manually_corrected_species_name
name_string_cleaned	name_string_cleaned
gbif_query_species_name	gbif_query_species_name
name_parser_timestamp	name_parser_timestamp
gbif_url	gbif_url
synonym_usage_key	synonym_usage_key
match_type	match_type
resolved_via_synonym	resolved_via_synonym
alternate_match_used	alternate_match_used
confidence	confidence
sim_name_confidence	sim_name
sim_authorship_confidence	sim_authorship
sim_classification_confidence	sim_classification
sim_rank_confidence	sim_rank
sim_status_confidence	sim_status

Table 6 - contents of dbo.predicts_2.name_resolution

The newly-minted name_resolution_id is returned when a row is written to name_resolution.

This key is used to update a third, occurrence-level table: species_record. Within this table, original_species_name is used to identify the appropriate line(s) to be updated:

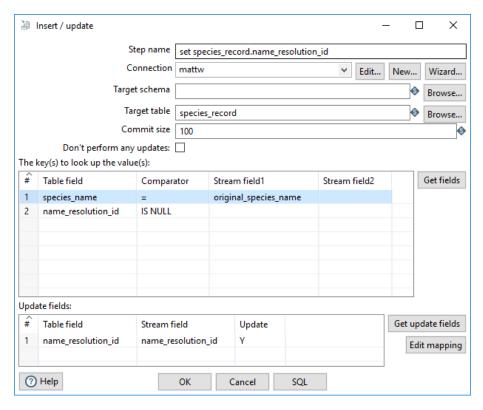


Figure 1 - Insert/update criteria for name_resolution

7.2 UNRESOLVED NAMES

Unresolved names fall into one of the following categories:

- a. No results: GBIF was unable to find a single match for this name
- b. Tied results: > 1 record has the highest confidence match value
- c. No reliable results: GBIF matches with a confidence score >= 94 (max. 100)

These are written to a new copy of needs_manual_name_resolution.xlsx_ for investigation + resolution by the P2 team. The fields included are:

field	description	
source_name	Name of research study/source from dbo.predicts_2.source	
non_standard_name_note	Indicator of type of name from dbo.predicts_2.species_record. E.g., 'Common name' or 'Uncertain taxa'	
high_level_taxa	Higher taxonomy indicator (can vary in level). Value from dbo.predicts_2.species_record	
original_species_name	Original species name from dbo.predicts_2.species_record	
manually_corrected_species_name	Null (to be populated by PREDICTS2 team)	
name_parser_timestamp	Date the pipeline was last run	
gbif_best_match	Match(es) with highest confidence suggested by GBIF. Ties results are '/' delimited.	
gbif_best_match_confidence	Confidence score of highest match suggested by GBIF	

Table 7 - composition of manual resolution spreadsheet