# *soweego*
# Link Wikidata to large catalogs

*Anywhere, 14 September 2021*

*Marco Fossati / Hjfocs*
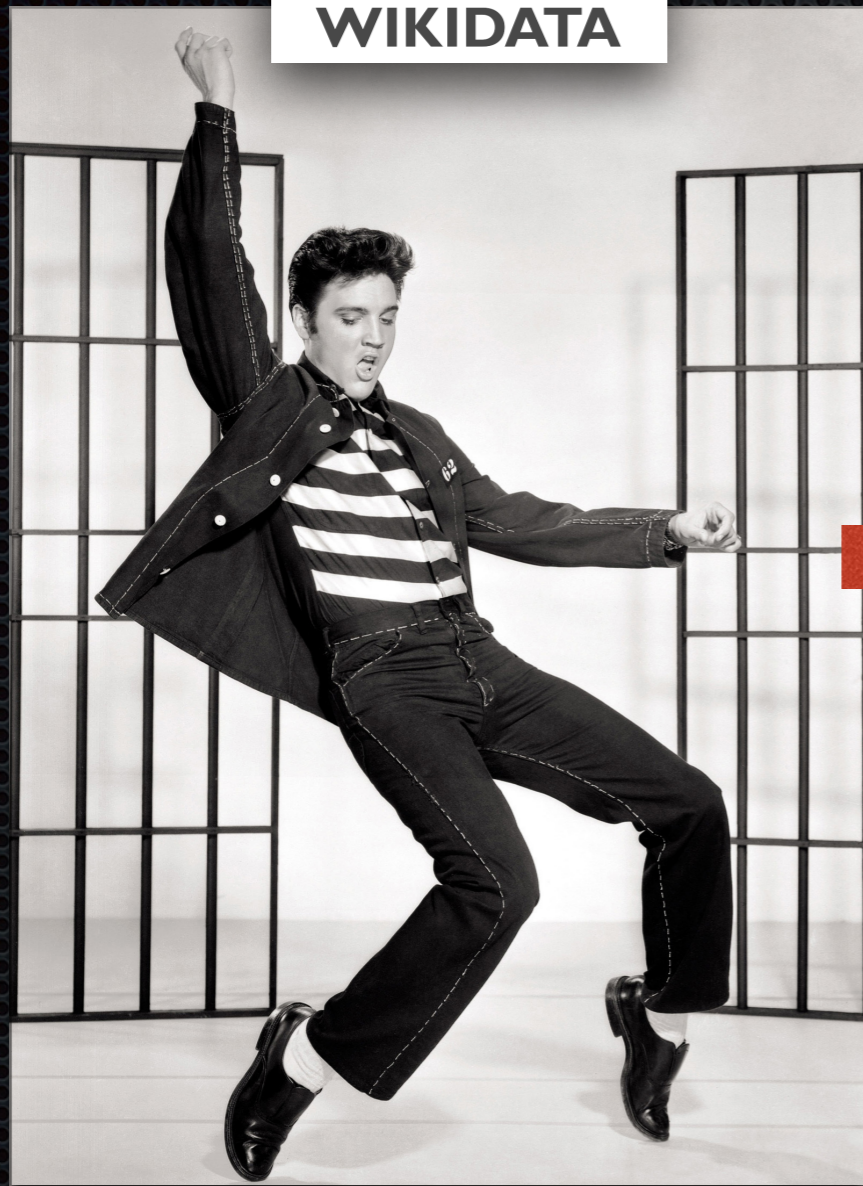
# Disclaimer

A tool or not a tool?

a WMF project
Project Grants program

# Example



WIKIDATA

MusicBrainz

Q303 → 01809

# Idea

- *Record linkage* AKA *data matching* AKA *entity resolution*

- Input: dataset **pair**

  - source = Wikidata

  - target = external catalog

- Output: **links**, as Wikidata identifier statements

# Example

# Why

- Increase Wikidata **quality** & **trust**

  - **quality** = identifiers enable feedback loops

  - **trust** = references to external reliable sources

# How

- Record linkage workflow

- Supervised machine learning

# What

- **People**, a big Wikidata slice

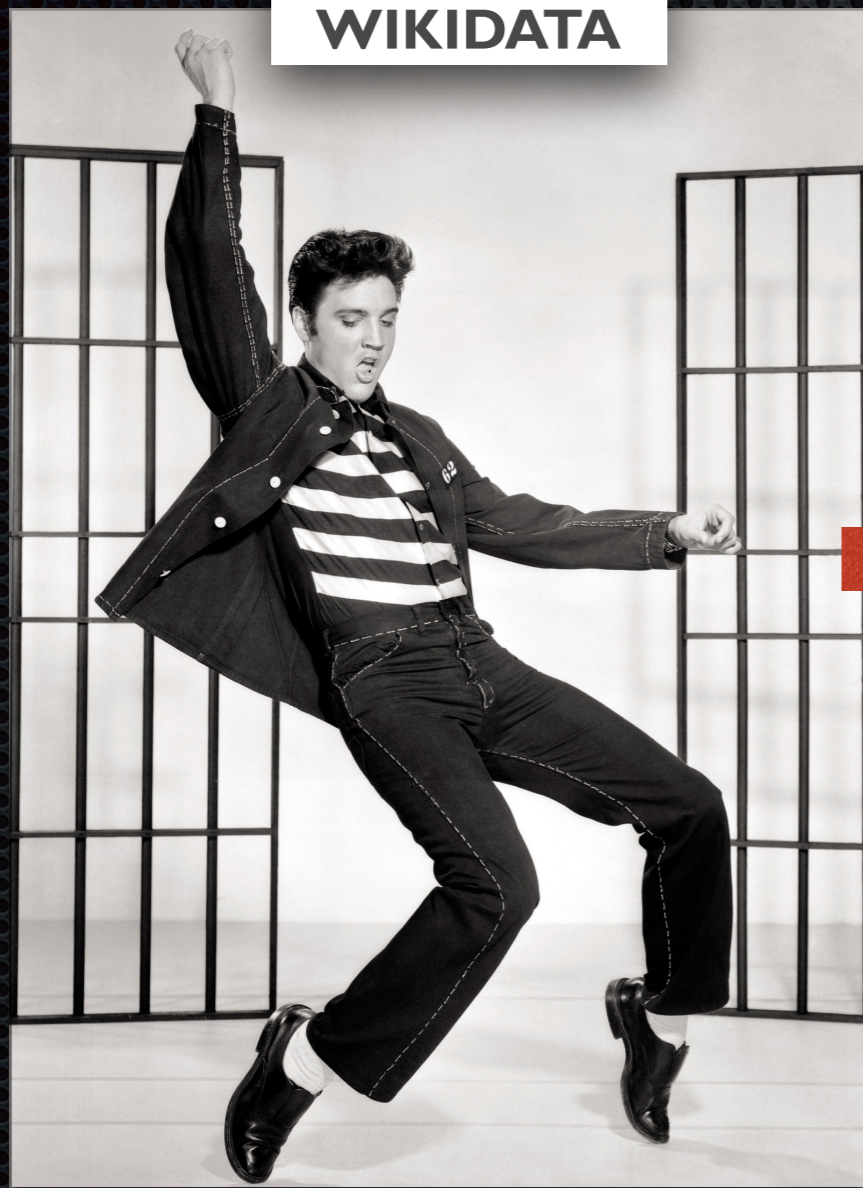- **Works**, made by people

# What's hot?

**Validator**: sync Wikidata to target catalogs

# In a nutshell

* **Run** validation criteria

  1. ID existence

  2. URLs

  3. domain-specific data

* **Apply** actions

  * enrich items

  * rank statements

  * submit to target communities

# Example



WIKIDATA

MusicBrainz

Q303 ➝ 01809

# #1 Existence

**Q303** > **01809**, *but* **01809** is *no more* in MusicBrainz

**WIKIDATA**

**Deprecate** statement

# #2 URLs

**Q303** has **7** URLs - **01809** has **8** URLs - *3* overlap

WIKIDATA

MusicBrainz

**add 5** URLs to **Q303**
(8 - 3)

**submit 4** URLs
(7 - 3)

# URLs are dangerous:
# let's add **IDs** first

# #3 Domain-specific data

**Q303** was born on **1935** in **Tupelo**
**01809** was born on **1934** in **Memphis**

**add 2** statements to
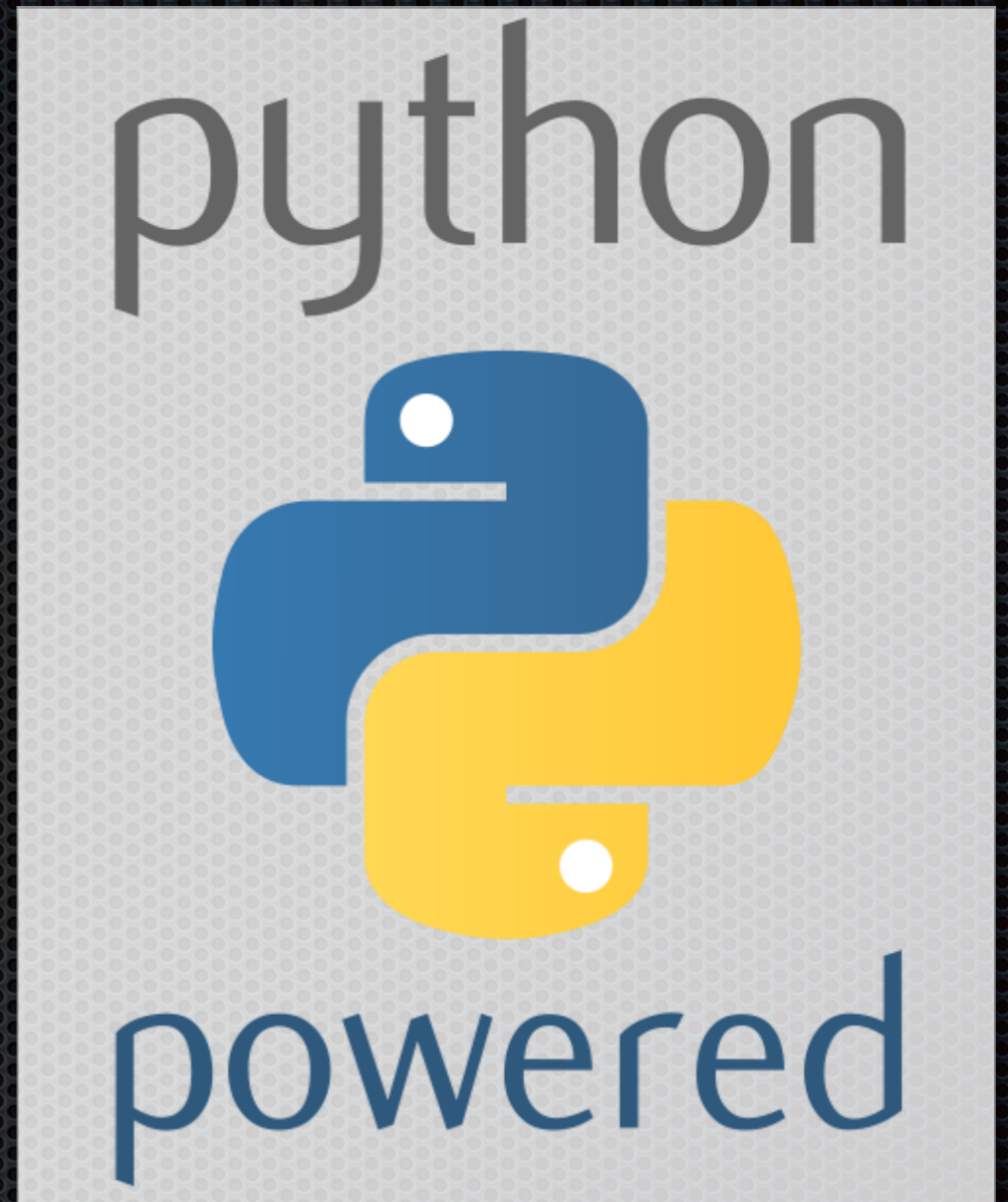**Q303**

**submit 2** values

# They can be controversial: let's feed the **mismatch finder** first

**WIKIDATA**

## Mismatch Finder

**About this tool**

The Mismatch Finder shows you data in Wikidata that differs from the data in another database, catalog or website (for example, someone's date of birth in Wikidata doesn't match the corresponding entry in the German National Library's catalogue). Mismatches like this need fixing, and the Mismatch Finder helps you to do just that.

# Thanks!

## *https://soweego.readthedocs.io/*

*Marco Fossati / Hjfocs*