# VanderBot:

## A spreadsheet-based system for creating and updating items in Wikidata

Steve Baskauf

steve.baskauf@vanderbilt.edu

Wikimedia username: Baskaufs

DISC DIGITAL SCHOLARSHIP AND COMMUNICATIONS

Jean & Alexander Heard LIBRARIES

VANDERBILT UNIVERSITY®

# Parts of the presentation:

1. Background: motivation, definitions, issues

2. Loading data into Wikidata with VanderBot: theory, workflow, and generalization

3. Maintaining an "authoritative" dataset in Wikidata

# This may be the basis for future blog posts

http://baskauf.blogspot.com/

- Initial posts based on rudimentary knowledge:
  - Python and SPARQL queries http://baskauf.blogspot.com/2019/05/getting-data-out-of-wikidata-using.html
  - Python and the Wikidata API http://baskauf.blogspot.com/2019/06/putting-data-into-wikidata-using.html
  - The "how can I use this?" question …
- Series documenting VanderBot, starting with:
  - http://baskauf.blogspot.com/2020/02/vanderbot-python-script-for-writing-to.html
  - and 3 subsequent posts
- Also see VanderBot repo landing page: http://vanderbi.lt/vanderbot

# Background:
## motivation, definitions, issues

# Why are people interested in Wikidata?

- Easy-to-use user interface for community contributions
- Strong community and high profile in knowledge graph community
- Interactivity with Wikidata and Wikimedia Commons
- "identifier central" (important for this group!)
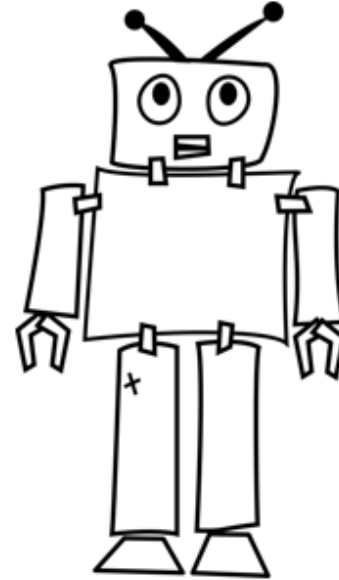- Linked Open Data (LOD) capabilities (e.g. SPARQL queries)

# Grab-bag of issues related to the VanderBot approach

- Do your data belong in Wikidata?
  - notability
  - scalability (vs. Wikibase)
  - linkability (would somebody ever link something to this?)
- Scalability of CSV2RDF (limits maybe 1M triples, 10 000 row CSV)
- Establishing the limits of the graph model (to fit "flat" spreadsheet)
- Practical matters of building a schema for your data consistent with the approach
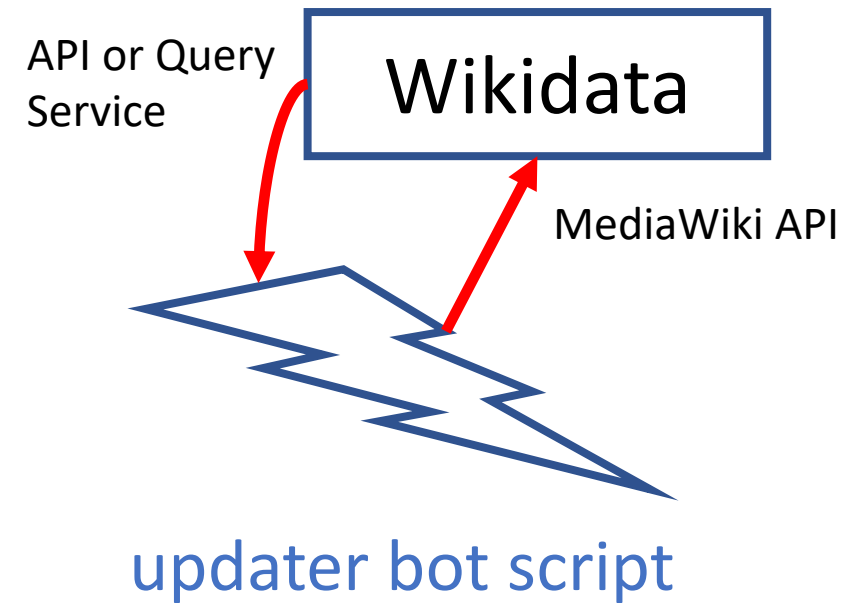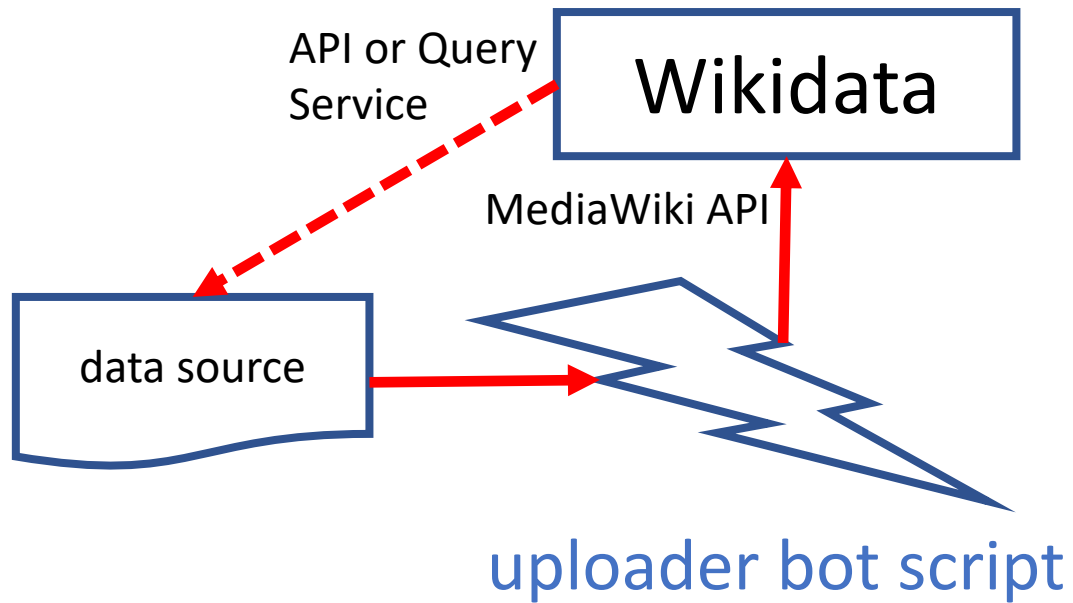- Importance of references

# "Authoritative" data into Wikidata

- Need to decide on data model/properties to be managed
    - Existing WikiProject with properties?
    - Query for properties of same type (e.g. what properties are used for coins?)
- Evaluate the current state of data in Wikidata
    - What fraction of possible items are already there?
    - Do all properties needed for the model exist?
    - Is there a value of a property (or property chain) that will link all items of interest?
- Evaluate "authoritative" data source for ease of linking
    - clean data
    - standardized strings
    - unique identifiers

# What's a bot?

- Software that can read and write to Wikidata via an API (application programming interface)
  - autonomous (no human intervention)
  - non-autonomous (human monitors and may intervene)

# "dumb" (autonomous) bot overview



API or Query Service

Wikidata

MediaWiki API

data source

uploader bot script

API or Query Service

Wikidata

MediaWiki API

updater bot script

# QuickStatements



Wikidata

web interface

human

keyboard or
web interface

MediaWiki API

CSV or other

QuickStatements
via web interface

OpenRefine is similar.

# VanderBot-human team overview

# Human editors



Greg Weldy (our star editor)
averages **160 edits per day**

high quality,
low speed

# "dumb" bot



Wikidata bot without a bot flag:
max **3000 edits per hour**

low quality,
high speed

# bot-human team



VanderBot

Create/edit about
**200 items per hour**
**(vs. 40 edits)**

optimizes quality
and speed

# Origins of the project



- Thinking about how Wikidata data could be archived locally (original focus on Wikibase)

- Previous (nonstandard) attempts to turn CSVs into RDF (Guid-O-Matic https://github.com/baskaufs/guid-o-matic)

- Inability to understand and use Pywikibot and Wikidataintegrator

- Frustration with labor-intensiveness and inability to script OpenRefine and QuickStatements



the Guid-O-Matic squid

# W3C Recommendation: Generating RDF from Tabular Data on the Web

- International standard to relate CSV tables to RDF

- https://www.w3.org/TR/csv2rdf/

- Approach: use a JSON-LD schema to map the table columns to the Wikibase graph model

- Applications:
  - VanderBot API-writing script knows how to convert CSV data to Wikibase-based JSON required by Wikidata API
  - Applying the schema to the CSV data will emit exactly the same RDF as is queried by the Wikidata Query Service

- Advantage: stable, standard way to unambiguously archive data in Wikidata in an easy-to-read tabular form, implementation independent.

# Generating a knowledge graph (RDF) from a CSV



spreadsheet (actively maintained or archived)

rdf-tabulator script

API-writing script

local triplestore (e.g. Fuseki)

Wikidata (now)

compare by federated SPARQL query

# Advantages

- Humans can review the spreadsheet to get a quick overview of the state of things (absence of references, frequency and diversity of values, label and description quality)

- Humans can change the spreadsheet to make many changes with few button clicks (copy and paste the same reference into hundreds of cells, quick cleanup of names and descriptions)

- The spreadsheet (along with mapping schema) is a snapshot of part of the Wikidata graph at a moment in time; can be reconstituted as RDF in a triplestore.

- Using Python scripts (vs. OpenRefine, Quickstatements) to have the potential to automate steps more fully.

# Example mapping



**JSON mapping file**

**CSC data file**

The JSON mapping schema makes it possible to "understand" what the table means in terms of the Wikibase model (specifically Wikidata).

**emitted RDF/Turtle (same as from Wikidata Query Service)**

# Generating the table schema



Using a web GUI (mostly unrestricted)

From simplified bespoke JSON (more restricted)

# VanderBot workflow (researchers)

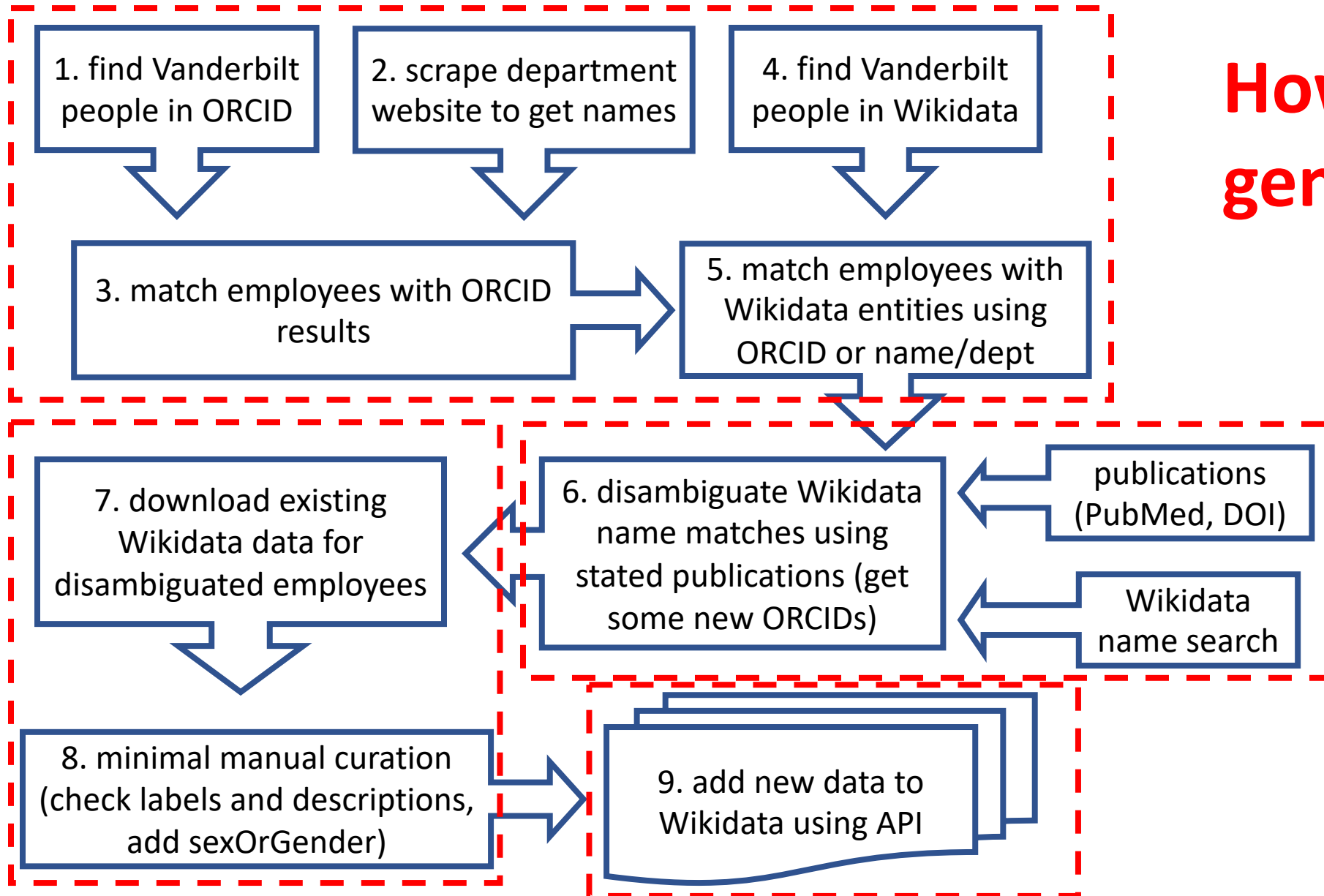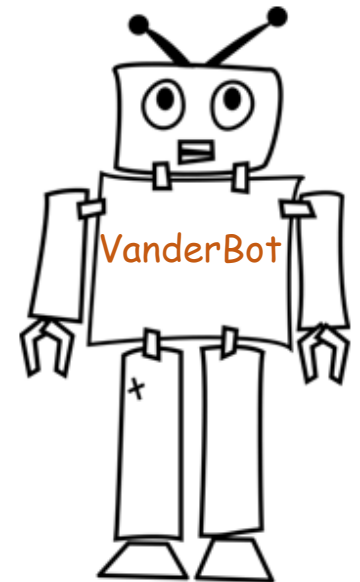- Over 8000 item edits (most editing multiple statements at once)
- Over 4600 Vanderbilt scholars/researcher items created or linked

1. find Vanderbilt people in ORCID

2. scrape department website to get names

4. find Vanderbilt people in Wikidata

3. match employees with ORCID results

5. match employees with Wikidata entities using ORCID or name/dept

6. disambiguate Wikidata name matches using stated publications (get some new ORCIDs)

publications (PubMed, DOI)

Wikidata name search

7. download existing Wikidata data for disambiguated employees

8. minimal manual curation (check labels and descriptions, add sexOrGender)

9. add new data to Wikidata using API

VanderBot

# VanderBot workflow

**How can this be generalized?**

1. find Vanderbilt people in ORCID

2. scrape department website to get names

4. find Vanderbilt people in Wikidata

3. match employees with ORCID results

5. match employees with Wikidata entities using ORCID or name/dept

7. download existing Wikidata data for disambiguated employees

6. disambiguate Wikidata name matches using stated publications (get some new ORCIDs)

publications (PubMed, DOI)

Wikidata name search

8. minimal manual curation (check labels and descriptions, add sexOrGender)

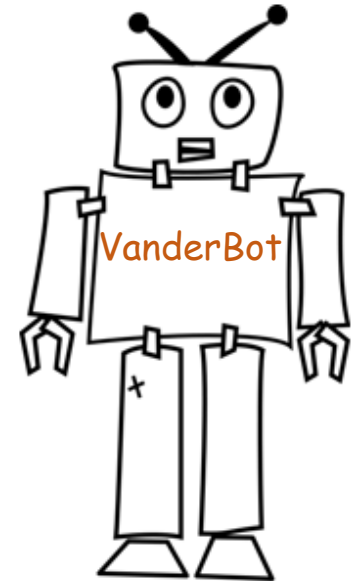9. add new data to Wikidata using API

VanderBot

# Generalized workflow

**How can this be generalized?**

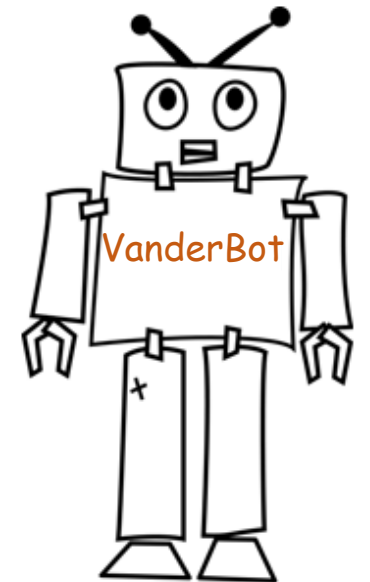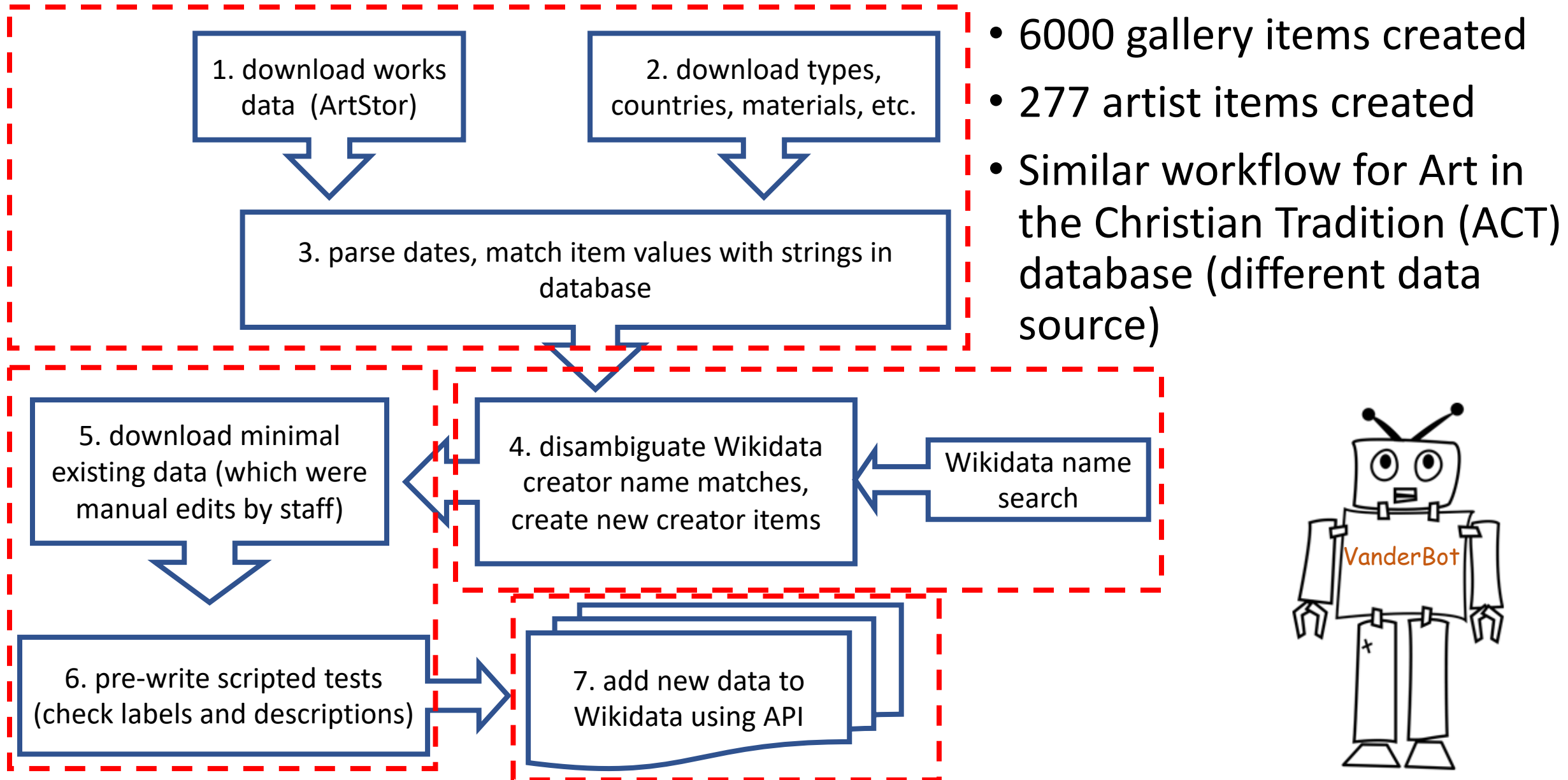Assemble authoritative source data (idiosyncratic)

Merge existing item data (SPARQL) with source data

Disambiguate against existing items

Add new data via API

VanderBot

# VanderBot workflow (Fine Arts Gallery)

1. download works data (ArtStor)

2. download types, countries, materials, etc.

3. parse dates, match item values with strings in database

5. download minimal existing data (which were manual edits by staff)

4. disambiguate Wikidata creator name matches, create new creator items

Wikidata name search

6. pre-write scripted tests (check labels and descriptions)

7. add new data to Wikidata using API

- 6000 gallery items created
- 277 artist items created
- Similar workflow for Art in the Christian Tradition (ACT) database (different data source)

VanderBot

# How can we keep track of existing items (of our "authoritative" interest)?

1. **Keep a list of Q IDs**
2. Use identifying property (or path) and value

Example:

```
SELECT DISTINCT ?qid ?label
WHERE {
VALUES ?qid
    {
    wd:Q102305506
    wd:Q102315563
    wd:Q102315787
    wd:Q102949359
    }
?qid rdfs:label ?label.
}
```

# How can we keep track of existing items?

1. Keep a list of Q IDs
2. **Use identifying property (or path) and value in triple pattern**

Examples:

- `?work wdt:P195 wd:`Q18563658. (collection Vanderbilt U. Fine Arts Gallery)
- `?researcher wdt:P1416/wdt:P749 wd:`Q29052. (affiliated with dept., parent organization Vanderbilt)
- `?image wdt:`P9092 `?id.` (Art in the Christian Tradition identifier)

# How to disambiguate strings against existing items?

- Generate a list of variant string forms (e.g. initials in names)

- Screen Wikidata SPARQL hits against criteria (human, not dead, born after…)

- Check hits against linked sources (PubMed, CrossRef, ULAN)

- Compare hits against existing items using fuzzy string matching.
  - Robust against small differences (periods, capitalization, diacritics)
  - Must test to determine best matching algorithm and cutoff score.
  - Not perfect (Bob vs. Robert, language variants)
  - Must accept the fact that mistakes will be made, but minimize.

- Present multiple possible matches for human decision.

- **Only partly generalizable, must experiment for each use case.**

Not unhuman?
(check "instance of")

Too old?
(check birth date)

Dead?
(check death date)

No bad description?
(e.g. "Ming Dynasty person)

# Software development

Assemble authoritative
source data (idiosyncratic)

Download existing
data
`acquire_wikidata_metadata.py`

Build csv-metadata
schema
`convert_json_to_metadata_sch`
`ema.py`

Disambiguate against existing
items
`vb3_match_wikidata.py`
`process_gallery.ipynb`

Add new data via API
`vb6_upload_wikidata.py`

VanderBot

# Maintaining an "authoritative" dataset in Wikidata

# Maintaining existing authoritative data

- "Authoritative data" is really an offense to the spirit of Wikidata, but...
- How do we detect changes from our authoritative data?
- How do we decide whether the changes are:
    - "good" (community contributions) that should be pulled into our local dataset?
    - "bad" (vandalism) that should be deleted/changed?
    - information outside the scope of our interest that should be ignored?
- How do we implement the transfer of new data?
    - Should a human intervene in every change?
    - What degree of automation is "safe"?
    - Should human approve changes in advance or review changes after the fact?

# Checking for changes over time



spreadsheet (archived)

rdf-tabulator script

addition

change

local triplestore (e.g. Fuseki)

Wikidata some time later

compare by federated SPARQL query

# Simple example (manual RDF handling)

1. Generate RDF from CSV and schema using **`rdf-tabulator`** Ruby script (written by Greg Kellogg) and save as an RDF/Turtle file.

2. Load the file into a Fuseki triplestore.

3. Generate the entailed "shortcut path" triples using SPARQL **`INSERT`** (loaded directly into triplestore by script).

4. Compare the local graph in Fuseki with the graph in Wikidata using a federated query to the Wikidata Query Service.

# 1. Generate RDF from CSV and schema

## Command line to run Ruby script:

```
rdf serialize --input-format tabular --output-format ttl --metadata csv-metadata.json --minimal > output.ttl
```



Turtle RDF serialization using Wikibase model

# 1. Generate RDF from CSV and schema

Conversion with `rdf-tabulator` is the rate-limiting step in process

| dataset | items | columns | CSV file size in kB (uncompressed/zip compressed) | conversion time (s) | triples | Turtle file size in kB (uncompressed/zip compressed) |
|---|---|---|---|---|---|---|
| Bluffton presidents | 10 | 32 | 9/3 | 1 | 193 | 24/4 |
| academic journals | 431 | 71 | 160/58 | 27 | 13 106 | 1500/320 |
| Vanderbilt researchers | 5247 | 30 | 3900/960 | 70 (1 min) | 62 634 | 7900/1300 |
| gallery objects | 4085 | 158 | 9300/1700 | 377 (6 min) | 209 597 | 28000/4400 |

Baskauf and Baskauf (in review) Table 2

Uninvestigated alternatives:

https://github.com/AtomGraph/CSV2RDF (151 348 939 triples in under 27 minutes)

https://github.com/Swirrl/csv2rdf (have not checked speed)

# 2. Load the file into a Fuseki triplestore

# 3. Generate the entailed "shortcut path" triples

```
with <http://researchers>
insert {?item ?truthyProp ?value.}
where {
    ?item ?p ?statement.
    ?statement ?ps ?value.
    filter(substr(str(?ps),1,40)="http://www.wikidata.org/prop/statement/P")
    bind(substr(str(?ps),40) as ?id)
    bind(iri(concat("http://www.wikidata.org/prop/direct/", ?id)) as ?truthyProp)
    }
```

wds:Q37371192-A0750BEF-303B-4665-8B34-80C1A3C39972
(statement instance)

p:P108
(property)

ps:P108
(property statement)

wd:Q37371192
(subject item)

wdt:P108
(direct property)

rdfs:label "Brandt F. Eichman"@en

wd:Q29052
(object item)

```
@prefix wd: <http://www.wikidata.org/entity/> .
@prefix wds: <http://www.wikidata.org/entity/statement/> .
@prefix wdt: <http://www.wikidata.org/prop/direct/> .
@prefix p: <http://www.wikidata.org/prop/> .
@prefix ps: <http://www.wikidata.org/prop/statement/> .
```
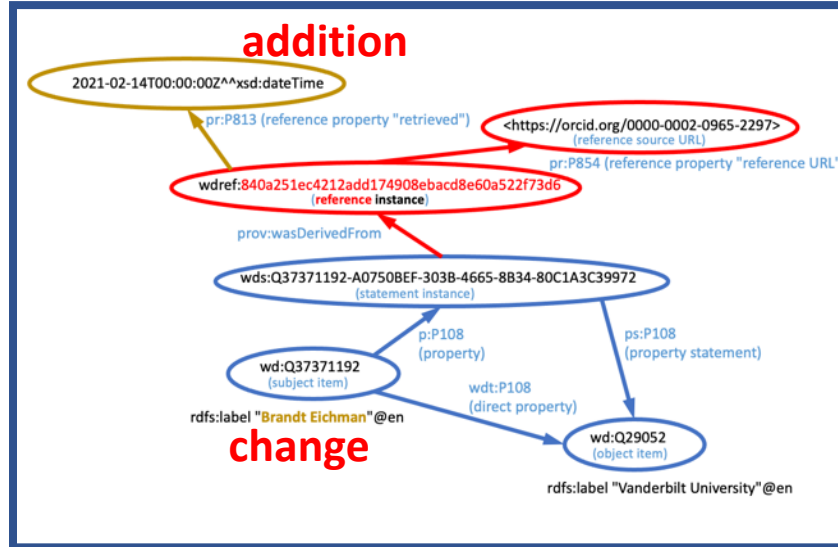
rdfs:label "Vanderbilt University"@en

# 4. Compare local graph with Wikidata graph

```
SELECT DISTINCT *
    WHERE {
    SERVICE <https://query.wikidata.org/sparql> {?subject ?predicate ?object.}
        MINUS
        {
        GRAPH <http://researchers> {?subject ?predicate ?object.}
        }
    }
```



Wikidata

graph generated from CSV

**Triples in Wikidata not in local copy (additions or changes)**
wdref:840a251ec4212add174908ebacd8e60a522f73d6 pr:P813 2021-02-14T00:00:00Z^^xsd:dateTime.
wd:Q37371192 rdfs:label "Brandt Eichman"@en.

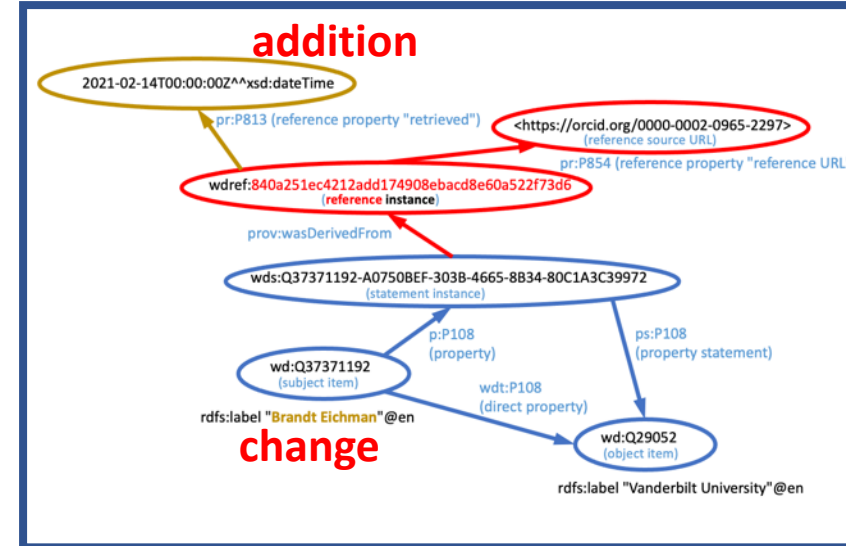# 4. Compare local graph with Wikidata graph

```
SELECT DISTINCT *
    WHERE {
    GRAPH <http://researchers> {?subject ?predicate ?object.}
      MINUS
      {
      SERVICE <https://query.wikidata.org/sparql> {?subject ?predicate ?object.}
      }
    }
```



graph generated from CSV

minus

Wikidata

**Triples in local copy not in Wikidata (deletions or changes)**

wd:Q37371192 rdfs:label "Brandt F. Eichman"@en.