# Applications
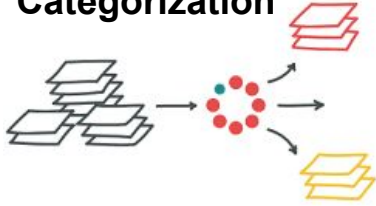
Amrapali Zaveri, Deniz Iren, Lea Beiermann

# Microtask Crowdsourcing

# Applications of Microtask Crowdsourcing

**Classification and Categorization**

**Finding Metadata**

**Ranking**

**Promoting**

**Data Collection and Enhancement**

**Sentiment Analysis**

Negative    Neutral    Positive
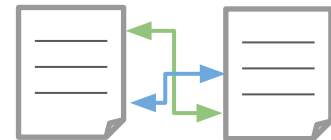
**Media Transcription**

**Content Feedback**

**Content Moderation**
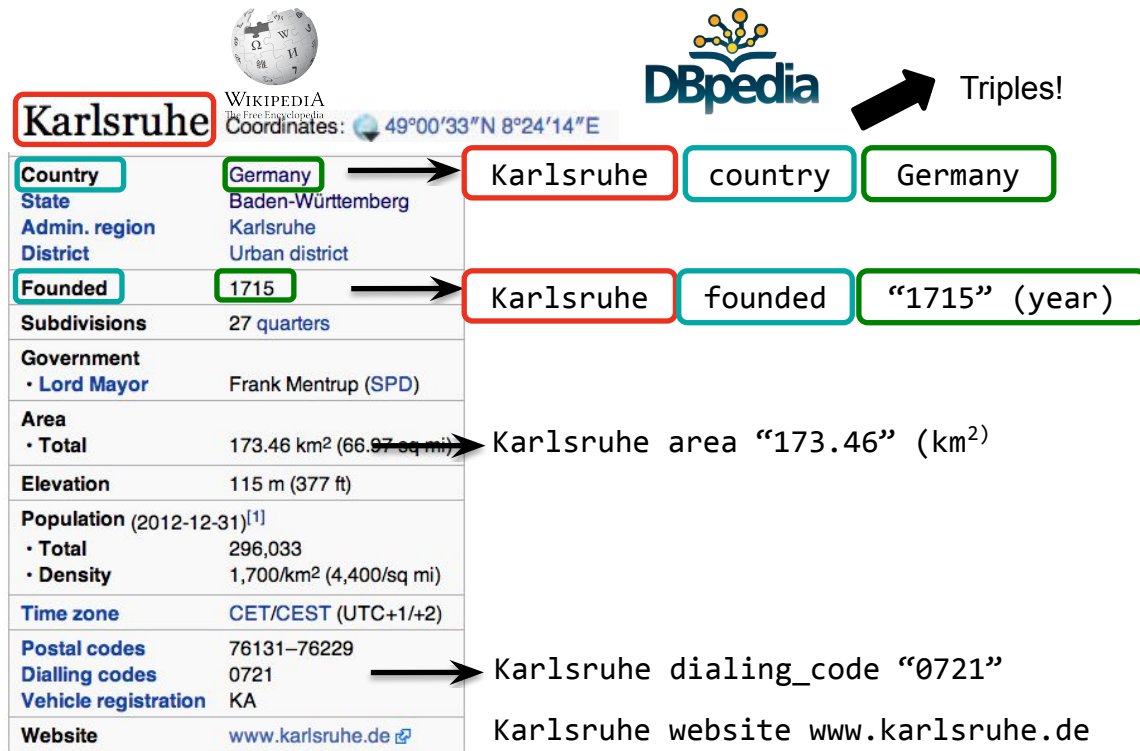
**Content Verification**

# Detecting Linked Data Quality Issues

# The DBpedia Knowledge Graph



http://en.wikipedia.org/wiki/Karlsruhe

Semi-structured data from Wikipedia

5

# The DBpedia Knowledge Graph



Triples!

| | |
|---|---|
| Karlsruhe | country | Germany |
| Karlsruhe | founded | "1715" (year) |

Karlsruhe area "173.46" (km$^2$)

Karlsruhe dialing_code "0721"

Karlsruhe website www.karlsruhe.de

# Quality Issues to Crowdsource

Three categories of quality problems occur
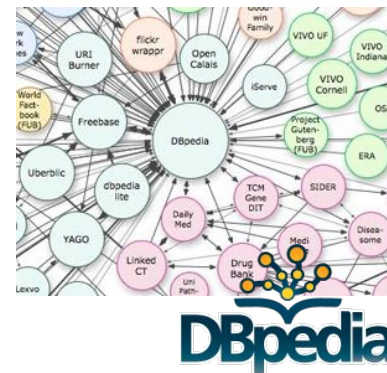in DBpedia [Zaveri2013] and can be crowdsourced:

- **Incorrect object**

  dbr:Dave_Dobbyn dbp:dateOfBirth **"3"** .

- **Incorrect data type or language tags**
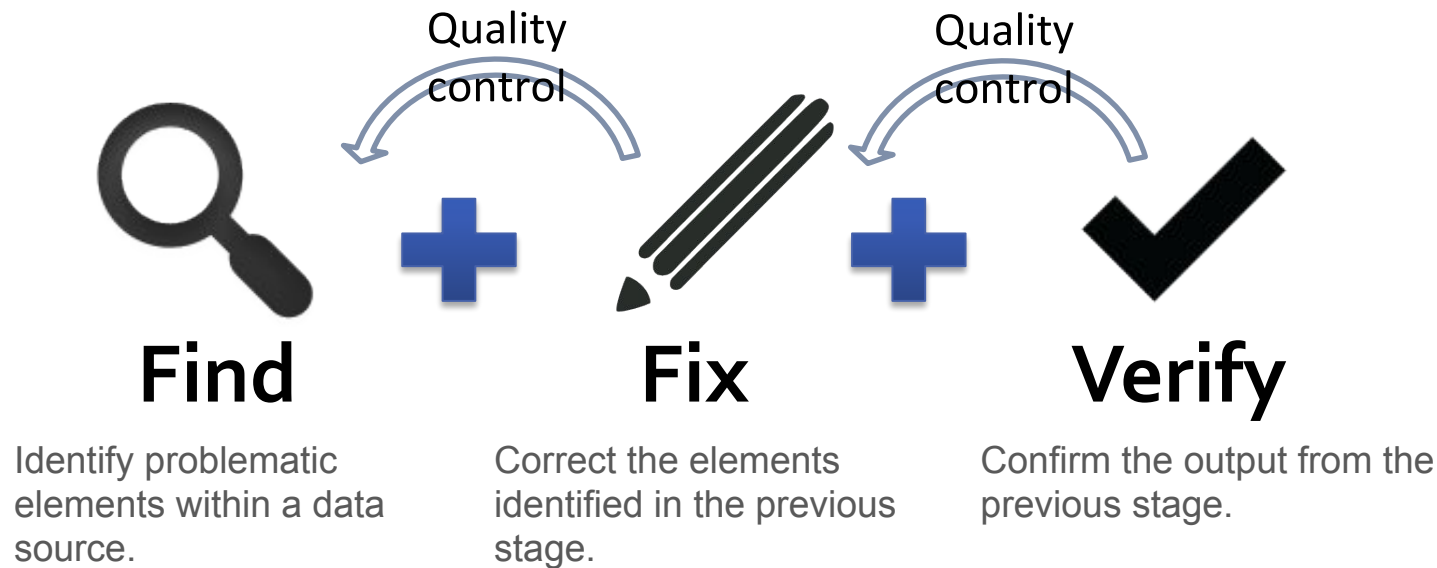
  dbr:Torishima_Izu_Islands foaf:name "鳥島"**@en** .

- **Incorrect link to "external Web pages"**

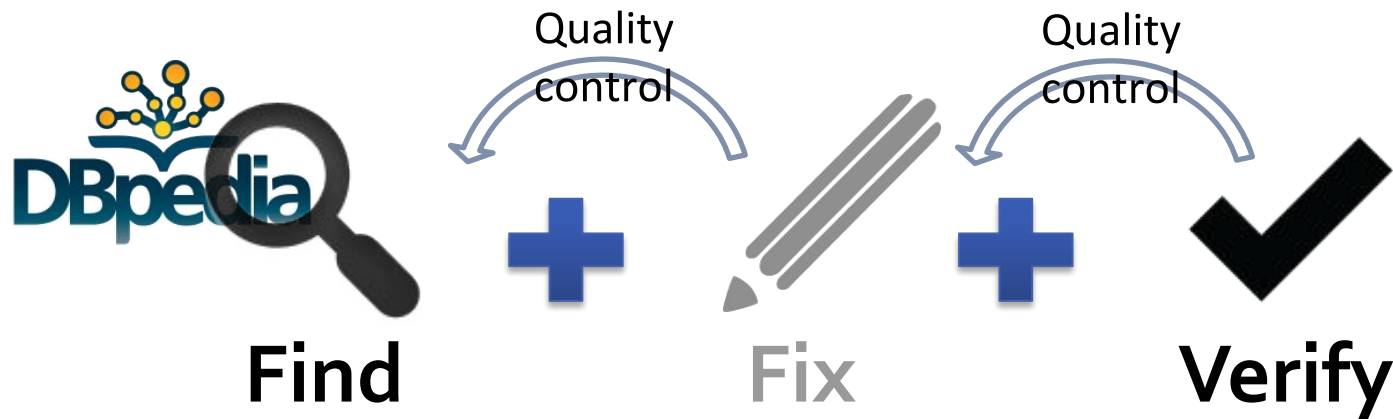  dbr:John-Two-Hawks dbo:wikiPageExternalLink **<http://cedarlakedvd.com>**.

# Crowdsourcing Approach

# Find-Fix-Verify Pattern [Bernstein2010]

Quality control

Quality control

## Find

## Fix

## Verify

Identify problematic elements within a data source.

Correct the elements identified in the previous stage.

Confirm the output from the previous stage.

# Applying Find-Fix-Verify to our Case Study: DBpedia-DQ

Quality control
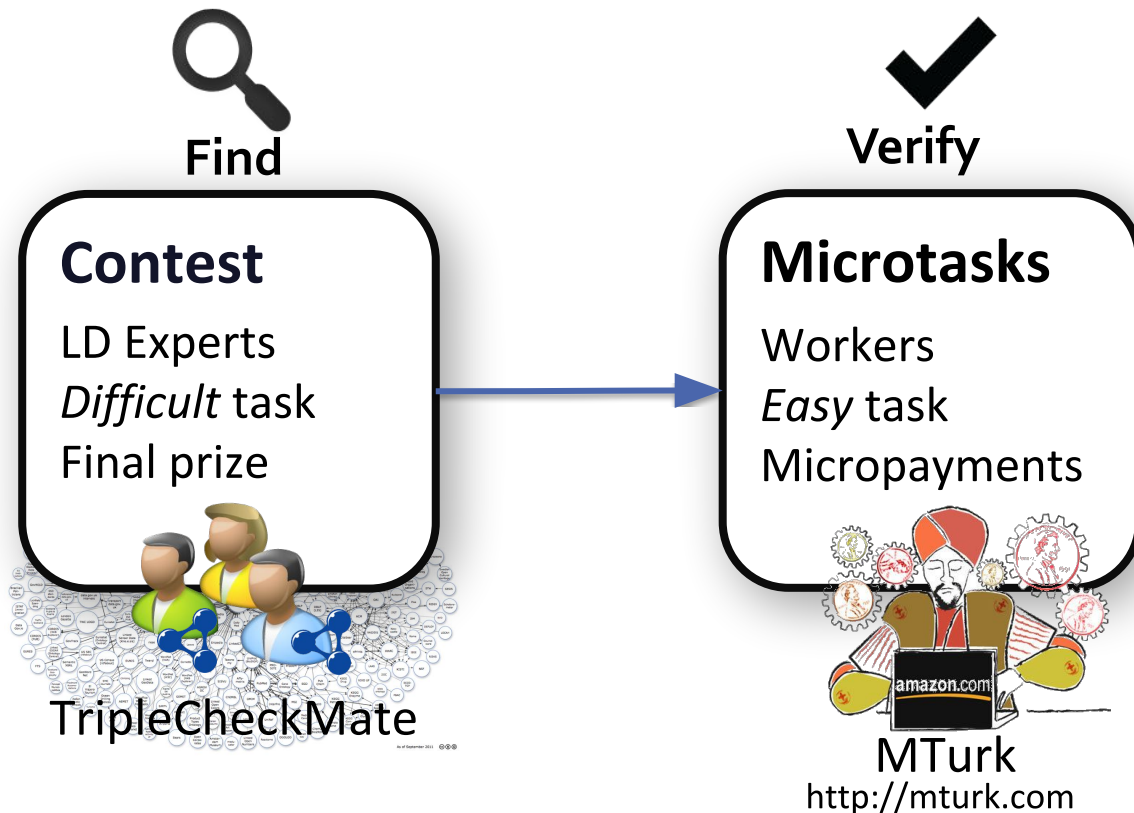
Quality control

**Find**

Identify erroneous triples and classify them according to the error found.

**Fix**

**Verify**

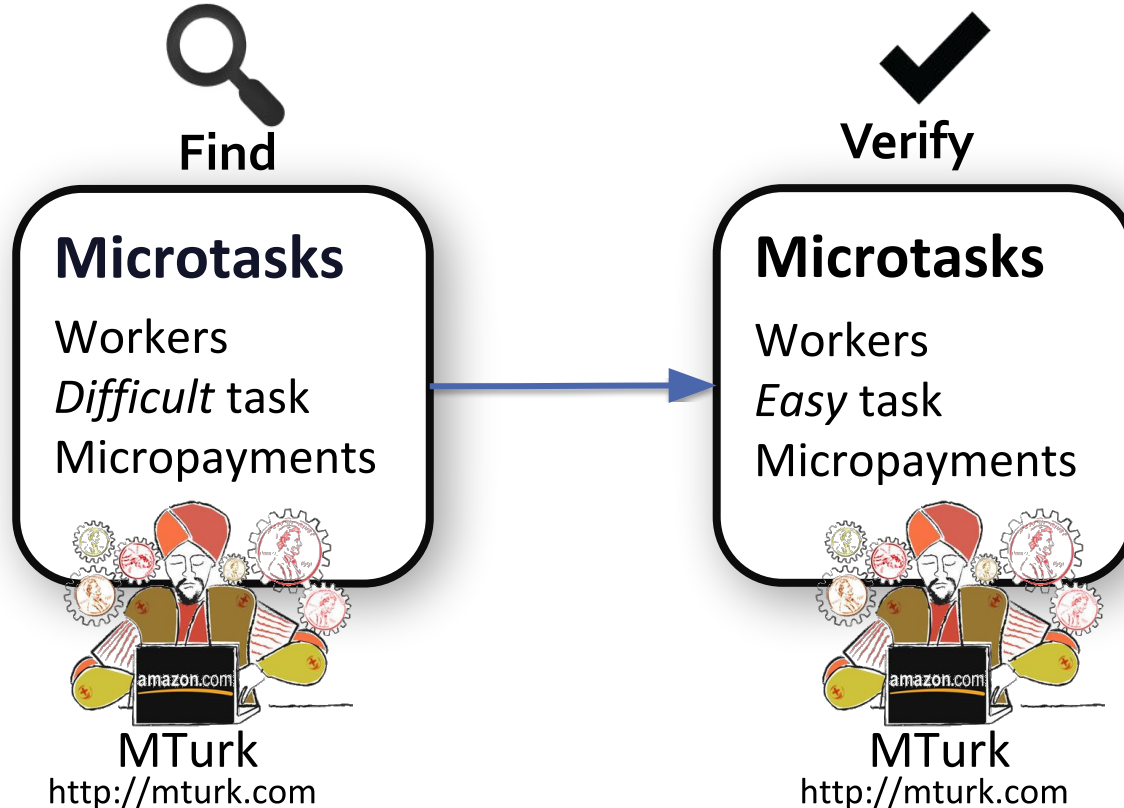Confirm the output from the previous stage.

DBpedia-DQ has two variants:
1. Combining experts + non-experts crowds (workers)
2. Using non-experts crowds (workers) in both stages

# Combining Experts + Workers (EW)

**Find**

**Verify**

### Contest

LD Experts
*Difficult* task
Final prize

TripleCheckMate

### Microtasks

Workers
*Easy* task
Micropayments

MTurk
http://mturk.com

# Combining Workers + Workers (WW)

**Find**

**Verify**

### Microtasks

Workers
*Difficult* task
Micropayments

MTurk
http://mturk.com

### Microtasks

Workers
*Easy* task
Micropayments

MTurk
http://mturk.com

# DBpedia-DQ Microtask Interfaces

**Find stage with workers:** MTurk Tasks

About: **Alexandria**

GO TO WIKIPEDIA ARTICLE: Alexandria

| WIKIPEDIA The Free Encyclopedia | DBpedia | Type of Errors |
|---|---|---|
| **Mar record low C:** *Not specified* | **Mar record low C:** 2 <br> Data type: Integer | ☐Value ☐Data type ☐Link |
| **Dec record high C:** *Not specified* | **Dec record high C:** 29 <br> Data type: Integer | ☐Value ☐Data type ☐Link |
| **Nov record low C:** *Not specified* | **Nov record low C:** 1 <br> Data type: Integer | ☐Value ☐Data type ☐Link |
| **Mar rain days:** *Not specified* | **Mar rain days:** 6 <br> Data type: Integer | ☐Value ☐Data type ☐Link |
| **single line:** *Not specified* | **single line:** yes <br> Data type: English | ☐Value ☐Data type ☐Link |
| **Aug record low C:** *Not specified* | **Aug record low C:** 18 <br> Data type: Integer | ☐Value ☐Data type ☐Link |

# DBpedia-DQ Microtask Interfaces

**Verify stage with workers:** MTurk Tasks

```
dbr:Dave_Dobbyn dbp:dateOfBirth "3" .
```

```
dbr:Torishima_Izu_Islands foaf:name "鳥島"@en .
```

```
dbr:John-Two-Hawks dbo:wikiPageExternalLink
    <http://cedarlakedvd.com>.
```

Incorrect object



Incorrect data type



Incorrect outlink



13

# DBpedia-DQ: Experimental Results



**Main findings:**
- It is difficult for the workers to assess datatypes
- Experts are not good in assessing external links
- Two-step validation increases the overall quality

**Main findings:**
- It is difficult for the workers to execute the find stage
- Workers are exceptionally good at identifying incorrect triples (high sensitivity)

M. Acosta, A. Zaveri, E. Simperl, D. Kontokostas, F. Flöck, J. Lehmann. Detecting Linked Data Quality Issues via Crowdsourcing. Semantic Web Journal, 2018.

# Experimental Results:
## Crowd-based vs. Automatic Data Quality Assessment

**Main finding:**



Humans (experts and workers) detected quality issues

that were not detected via RDFUnit (automatic tool)

and vice versa.

M. Acosta, A. Zaveri, E. Simperl, D. Kontokostas, F. Flöck, J. Lehmann. Detecting Linked Data Quality Issues via Crowdsourcing.Semantic Web Journal, 2018.

# Deniz

# Feedback Please !

http://bit.ly/crowdsourcing-feedback