# Knowledge Base Completion from Web Tables
## An experiment with the Yahoo Knowledge Graph

DBpedia Meetup 10/27/2016
Work by Dominique Ritze - Presented by Nicolas Torzec

# About Us...

**Dominique Ritze**

- Ph.D. student, **University of Mannheim, Germany**
- Knowledge base completion from Web tables: **T2K**
- Internship on Slot Filling in the Yahoo Knowledge Graph team

**Nicolas Torzec**

- Lead the Data / Science team for the **Yahoo Knowledge Graph**
- Information extraction, knowledge graph, semantic search, etc.
- NLP+TTS⇒Shopping Comparison⇒Web Search⇒Knowledge Graph

# Yahoo Knowledge Graph

**Large unified knowledge base**
- millions of entities, billions of facts
- power knowledge-based services at Yahoo

**Knowledge Acquisition**
- extract info about entities from feeds/web
- type inference, schema matching, normalization

**Knowledge Integration**
- cluster entities that refer to same real-world entity
- blend/fuse entity clusters: their types and facts

**Knowledge Enrichments**
- notableFor, related entities, hero images
- editorial curation!

**Serving...**

# Knowledge Base Completion from Web Tables

## Goal: enrich KG with information extracted from tables on the Web

- i.e. some information are only (or most easily) available in Web tables.
- **Slot filling** | Set expansion | Attribute discovery | Source discovery | etc.



```
<!DOCTYPE html>
<html>
<table>

</table>
</html>
```

# Knowledge Base Completion from Web Tables (2)

| Airport Name | Country Name | City | World Area Code | Airport Code |
|---|---|---|---|---|
| Puerto Rico Airport | Colombia | Puerto Rico | 327 | PCC |
| Uribe Airport | Colombia | Uribe | 327 | URI |
| El Nuevo Dorado International Airport | Colombia | Bogota | 327 | BOG |

http://www.altiusdirectory.com/Travel/colombia-airports.php

| City served / Location | Department | ICAO | IATA | Airport name |
|---|---|---|---|---|
| Manizales | Caldas | SKMZ | MZL | La Nubia Airport |
| Uribe | Meta | SKUB | URI | Uribe Airport |
| Corozal | Sucre | SKCZ | CZU | Las Brujas Airport |

https://en.wikipedia.org/wiki/List_of_airports_in_Colombia

YAHOO! KNOWLEDGE    PRODUCTION ▾

```
▼ object {19}
    @id : XXXXXXXXX
    label : Uribe Airport
  ▶ alternateEntityName [2]
  ▶ @type [5]
  ▶ source [2]
  ▶ alternateKey [2]
    description : Uribe Airport (URI) is an airport located in Uribe, Meta, Colombia.
    latitude : 3.21667
    longitude : -74.4
    eastLongitude : -74.39545
    westLongitude : -74.40455
    southLatitude : 3.21212
    northLatitude : 3.22122
    area : 1000000
  ▶ timeZone [1]
  ▶ countyLocation [1]
  ▶ countryLocation [1]
  ▶ cityLocation [1]
  ▶ stateLocation [1]
```

No IATA code for this airport...

YAHOO!

# Overview



**Slotting**

**Knowledge Base**

**Extraction**

Extract relational tables

(from Web crawl)

**Matching**

Match table rows/columns
to KG entities/properties

(type|schema|instance matching)

**Fusion**

Deduplicate and identify most likely
value of each table entity/property

(across tables)

# Table Extraction

## Goal: extract <u>relational</u> tables

- e.g. layout tables vs. relational tables vs. single-entity tables vs. ...



## Classification problem

- i.e. distinguish relational tables from other tables
- Features: level of nesting, avg cell length, % of links, datatype heterogeneity, etc.
- 11B tables in WDC Common Crawl 2012 ⇒ **1.3% are relational** (~91M tables)

**Slotting**

**Knowledge Base**

**Extraction** → **Matching** → **Fusion**

YAHOO!

# Matching

**Goal = match table/rows/columns to type, entities, and properties**

| WOO:/Entity/Place/Structure/TransportHub/Airport | | | | |
|---|---|---|---|---|
| label | countryLocation | location | ? | iataCode |
| **Airport Name** | Country Name | **City** | **World Area Code** | **Airport Code** |
| Puerto Rico Airport | Colombia | Puerto Rico | 327 | PCC |
| Uribe Airport | Colombia | Uribe | 327 | URI |
| El Nuevo Dorado International Airport | Colombia | Bogota | 327 | BOG |

YK:ID1
YK:ID2
YK:ID3

# Matching: Preprocessing

## Table preprocessing

- Filter low-quality tables (domain, page, language, number of rows/cells)
- Standardize cells (datatype casting, unit conversion)
- Detect header row, key column (based on datatype, uniqueness, position)
- ⇒ **Subset of 30M tables** (589M rows / 2.6B cells)

| Airport Name | | Country Name | City | World Area Code | Airport Code |
|---|---|---|---|---|---|
| Puerto Rico Airport | | Colombia | Puerto Rico | 327 | PCC |

## KB preprocessing

- Filter entities based on provenance, language, type, richness
- Filter properties irrelevant for matching ; denormalize object properties
- ⇒ **Subset of 56M entities** (854M facts) (~10x DBpedia)

# Matching: Blocking (1)

## Identify most likely "KG entities" for each "table entity"

- Why: to prevent full pair-wise comparison (i.e. **589M rows x 56M entities)**
- How: exact match on label/aliases, or multi-attribute fuzzy matching



- ⇒ **130M rows x 21M entities ; median number of matches: 7** (vs. DBpedia: 1)

# Matching: Blocking (2)

## Identify most likely "KG types" for each table

- Why: to limit further the number of comparisons and improve precision.
- Features: types of matched entities, type specificity/relatedness, etc.

|  | Airport | Place | Person | MusicTrack |
|---|---|---|---|---|
| Puerto Rico Airport | 3 | 3 | 0 | 1 |
| Uribe Airport | 1 | 1 | 0 | 1 |
| El Nuevo Dorado International Airport | 1 | 1 | 0 | 0 |

## Based upon most likely types:
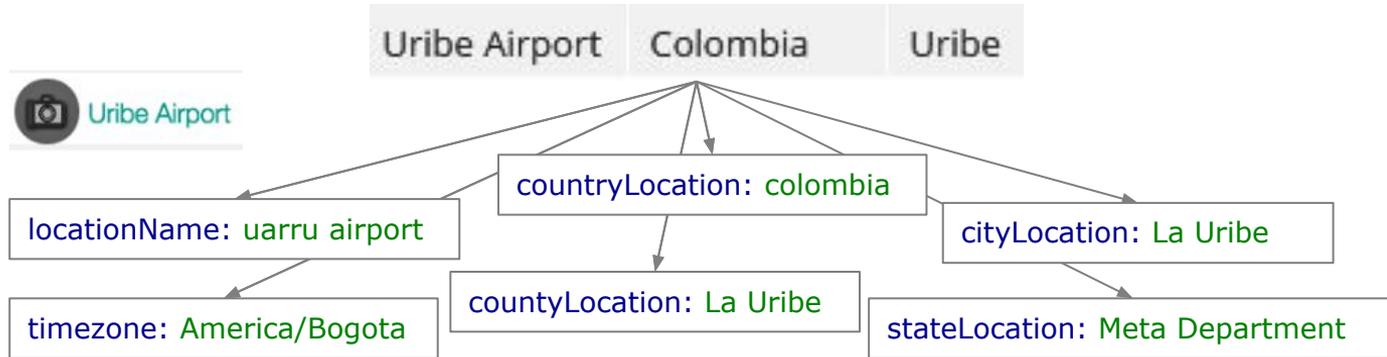
- Filter out rows whose type does not match the type of the table
- Redo entity blocking, focusing on entities whose type match one of the most likely types.

# Matching: Type Inference & Schema Matching (1)

**Assign a KG type to each table and a KG property to each column**

1) Compute similarities between table cells and entity facts
- Datatype-specific comparisons
- Set comparisons

# Matching: Type Inference & Schema Matching (2)

**Assign a KG type to each table and a KG property to each column**

2) Infer column-to-property mapping
- Naive approach: aggregate cell/fact similarity at the column level (⇒ ML)

3) infer table to type mapping
- Naive approach: aggregate column/property similarity at the table level (⇒ ML)

| Column: city | Puerto Rico Airport | | | Uribe Airport | El Nuevo Dorado Int. Airport | |
|---|---|---|---|---|---|---|
| **KG property** | **YK1** | **YK2** | **YK3** | **YK4** | **YK5** | **score** |
| cityLocation | 1.0 | 0.8 | 1.0 | 0.5 | 0.7 | 4.0 |
| countryLocation | 0.0 | 0.2 | 0.0 | 0.5 | 0.1 | 0.8 |

# Matching: Instance Matching

## Match each table entity to a KG entity

- Classification problem
- Features: property similarity score, value similarity score, entity type, etc.

| Country Name | City |
|---|---|

0.9     0.6

| countryLocation | cityLocation |
|---|---|

| Uribe Airport | Colombia | Uribe |
|---|---|---|

1.0     0.5

| countryLocation: colombia | cityLocation: La Uribe |
|---|---|

*naive similarity score ⇒ 0.9\*1.0 + 0.6\*0.5*

- **Naive classifier ⇒ ML classifier**

# Matching: Some Results

## Naive preprocessing / blocking
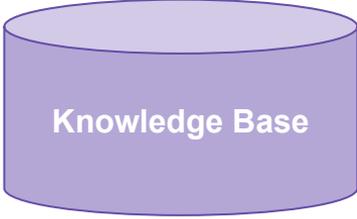
- Tables: from 91M tables to 30M tables (589M rows x 2.6B cells)
- KG: subset of 56M entities x 854M facts (~10x DBpedia)
- Blocking ⇒ 130M rows x 21M entities (median number of matches per row: 7)

## Naive matching

- table2type: 7M correspondences
- column2property: 9M correspondences
- row2entity: 82M correspondences

<span style="color:red">Room for improvement: new features, ML, joint optimization</span>

- Number of generated triples: 144M
- Number of unique value across entity-property group: 13M
- ~16% of the groups do not overlap with YK

<span style="color:red">Need data **fusion** to find most likely value for each group</span>

# Fusion

## Goal: identify most likely value for each table entity/property group

- Local Closed World Assumption (LCWA): everything in the KB is true
- Knowledge Base Trust (KBT) to assign a reliability score to a column
- Filter out not reliable columns
- "Most likely value" = weighted majority/median (with KBT score as weight)



- Works best for large tables and facts coming from many different tables!

# Fusion Some Results

## Fusion results (with naive approach)

- Metrics: precision/recall based upon overlapping facts
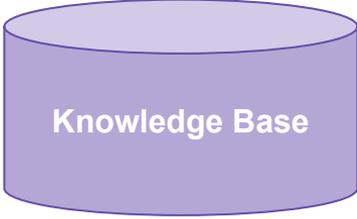- Slight deviations allowed...

| | Missing facts | Overlapping facts | Precision |
|---|---|---|---|
| MediaRelease.albumTrackOf | 65,867 | 105,484 | 0.92 |
| EducationalInstitution.countryLocation | 1336 | 1313 | 0.99 |
| CreativeWork.releaseDate | 0 | 208,626 | 0.74 |

## Expected slotting (with naive approach)

- 190K new/missing facts with exp. precision of 0.99 (2% of all facts)
- 480K new/missing facts with exp. precision of 0.86 (4% of all facts)

# Take Aways

- Large amounts of facts available in Web tables
- ⇒ Knowledge Base Completion from Web tables

- Naive end-to-end system running on Yahoo Knowledge Graph
- Slot filling: 480K new facts @ expected precision of 0.86
- Decent baseline but improvements are possible at each step!

- Future directions:
  - ML and KG as feedback loop
  - Joint optimization problem
  - Topic/Domain-specific applications

YAHOO!

# Thank You.

torzecn@yahoo-inc.com
twitter: @nicolastorzec

yes, we are hiring too.

YAHOO!