# User-Centric Ontology Population

KENNETH CLARKSON, **ANNA LISA GENTILE**,
DANIEL GRUHL, PETAR RISTOSKI,
JOSEPH TERDIMAN, STEVE WELCH

IBM RESEARCH

# Motivation

"The Semantic Web provides a common framework that allows data to be shared and reused across application, enterprise, and community boundaries" - Tim Berners-Lee 2001

17 years later still vast amount of valuable unstructured and semi-structured data is published on the Web

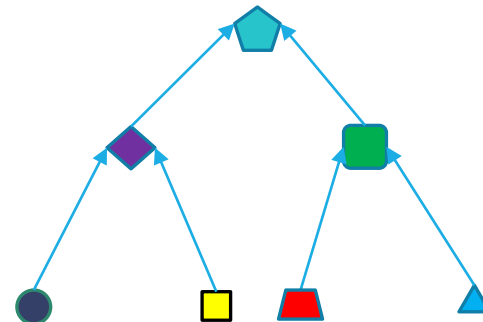**Goal**: automatically extract semantic data from text

# Problem Statement

# Problem Statement

### Entity Detection



### Target Ontology

# Problem Statement

User-defined Entity Grouping

Target Ontology

# Problem Statement

Ontology Alignment

Target Ontology

# User concepts

# Ontology Population

# Knowledge Discovery Process
(Fayyad et al. 1996)

# How good are machines?

~80% accuracy

~85% accuracy

~90% accuracy

# Is 80% enough?

# Introduce the human-in-the-loop

# Introduce the human-in-the-loop

"Computers are incredibly fast, accurate, and stupid.
Human beings are incredibly slow, inaccurate, and brilliant.
Together they are powerful beyond imagination."

Einstein never said that

# Vision

# Proposed Solution

Given initial user conceptualization, the methodology supports:
- Finding candidate ontologies
- **Aligning** the user's conceptualization to a target ontologies
  - novel <span style="color:red">hierarchical classification approach</span>
- **Maintenance** lifecycle
  - **build** (create new concepts)
  - **change** (splitting/merging concept)
  - **grow** (adding new instances to each concept)
    - from target ontologies
    - new facts extracted from unstructured data

# Implementation



Text corpus — Entity extractor — Extracted entities — User's conceptualization — Scout for ontologies — Select an ontology

Ent1
Ent2
.
.
EntN

1. Align knowledge

2. Maintain knowledge

# Aligning Input with a Target Ontology

Identify available ontologies
- collective instance matching

Align user conceptualization using ML models
- **Training data**: instances of the target ontology
- **Features** : domain-specific word embeddings
- **Classification strategies**
  - Flat hierarchical classification
  - Top-down local classifier per parent node
  - Combine flat hierarchical with top-down local classifier per parent node
- **Classifiers**
  - SVM, Random Forests, Logistic Regression, Convolutional Neural Network

# Flat Hierarchical Model

One model **for each level** of the hierarchy

- ◦ Simple
- ◦ High model complexity down the hierarchy
- ◦ precision declines

# Top-Down Local Classifiers

One model **for each parent** in the hierarchy

- ◦ Simple
- ◦ Error propagation through levels

# Combine Both Models

**Combine flat hierarchical models with top-down local classifier**

- flat model for level L-1
- local model for level L

# Ontology Maintenance

Adding new instances
- Use existing models

Reassigning Instances
- Leave-one-out validation

Generating new concepts
- If the class distribution is uniform then search for new concept

Merging concepts
- User's concepts aligned to the same target ontology concept should be merged

Concept splitting
- Use hierarchical clustering
- Refine until a criteria is met

# Evaluation - Alignment

Task: label adverse drug events with preferred medical terms

Data:
- MedDRA ontology as a target ontology
- ADE groups extracted from "ask a patient blogs"

|  | User's conceptualization | MedDRA |
|---|---|---|
| #level1 | 17 | 27 |
| #level2 | 62 | 304 |
| #level3 | 106 | 1,444 |
| #level4 | 169 | 20,935 |
| #Instances | 3,262 | 95,061 |

Evaluation metric: HITS@10 = proportion of correct mapping top 10 ranked suggestions
- Evaluate per each level of the hierarchy



**System Organ Class**
Gastrointestinal disorders

**High Level Group Term**
Gastrointestinal signs and symptoms

**High Level Term**
Nausea and vomiting symptoms

**Preferred Term**
Nausea

**Lowest Level Term**
Feeling queasy

# Evaluation – Ontology Alignment

Baselines:
- String-based average-link matching
- Word embeddings
- LDA topic modeling

Evaluation metric: HITS@10
- proportion of correct mappings that appear in the top 10 ranked suggestions
- Evaluate per each level of the hierarchy

# Results

# Results: Level 4

# Evaluation – Ontology Maintenance

**Adding new instances** – evaluate how precise the models can add new instances to the already aligned concepts

- Retrieved 298 **new** ADE from askapatient.com
- Measure HITS@k for each level of the hierarchy

# Evaluation – Ontology Maintenance

**Adding new concepts** – evaluate the model's ability to notify the user to add a new concept

- Evaluation
- Selected 500 MedDRA instances that don't belong to the user's conceptualization (positive instances), and 500 instances that belong to the user's conceptualization

Results:

- Precision: 73.8%
- Recall: 84.6%
- F-score: 78.83%

$$E(x) = \sum_{i=0} kP(C_1|x) * log_2 P(C_1|x) > 1$$

# Evaluation – Ontology Maintenance

**Adding new concepts**

- model's ability to suggest the user to add a new concept
- evaluation
  - 500 MedDRA instances that don't belong to the user's conceptualization (positive instances)
  - 500 instances that belong to the user's conceptualization

Results:

- Precision: 73.8%
- Recall: 84.6%
- F-score: 78.83%

# Evaluation – Ontology Maintenance

**Re-assigning Instances**: evaluate the model's ability to reassign instances to other concepts.

The model identified 82 instances to be reassigned, from which 67 (81.7%) were accepted by the medical doctor

Examples:

◦ User errors: "*stomach aches*" was assigned to "*Emotional disorder*", which should be assigned to "*Abdominal distension*"

◦ Better matches: "*sensitivity to light*" was assigned to "Visual impairment", which was later reassigned to "Photophobia"

# Further Use-Cases

- **Maintain health and medical data**
  - Adverse drug reactions
  - Drug brands
- **Maintain e-shop product catalog and taxonomy**
  - Map new features to an existing product catalog
  - Map new products in the product taxonomy
- **Social media analysis**
  - Identifying new trends
- **Reviews analysis**
  - Movies and actors

# Conclusion

User-centric ontology population

**Human-in-the-loop** for each step
- Building, connecting and maintaining their conceptualization, using available ontologies
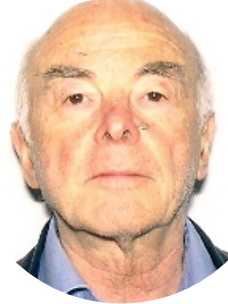
Novel **hierarchical classification model**
- dynamically refined based on user interaction

The approach supports the user to achieve **nearly perfect performance**

The user has full control on their level of involvement in the process
- Trade-off between involvement/cost/time and performance/quality of results

IBM

# User-Centric Ontology Population



Kenneth Clarkson, **Anna Lisa Gentile**, Daniel Gruhl,Petar Ristoski,Joseph Terdiman, Steve Welch

*annalisa.gentile@ibm.com*     *@AnLiGentile*