# Non-Parametric **Class Completeness** Estimators for Collaborative **Knowledge Graphs**

## The Case of Wikidata

**Michael Luggen,** Djellel Difallah, Cristina Sarasua,
Gianluca Demartini, and Philippe Cudré-Mauroux

ISWC 2019, Auckland

1

# Agenda

- Motivation

- Species Richness Estimators

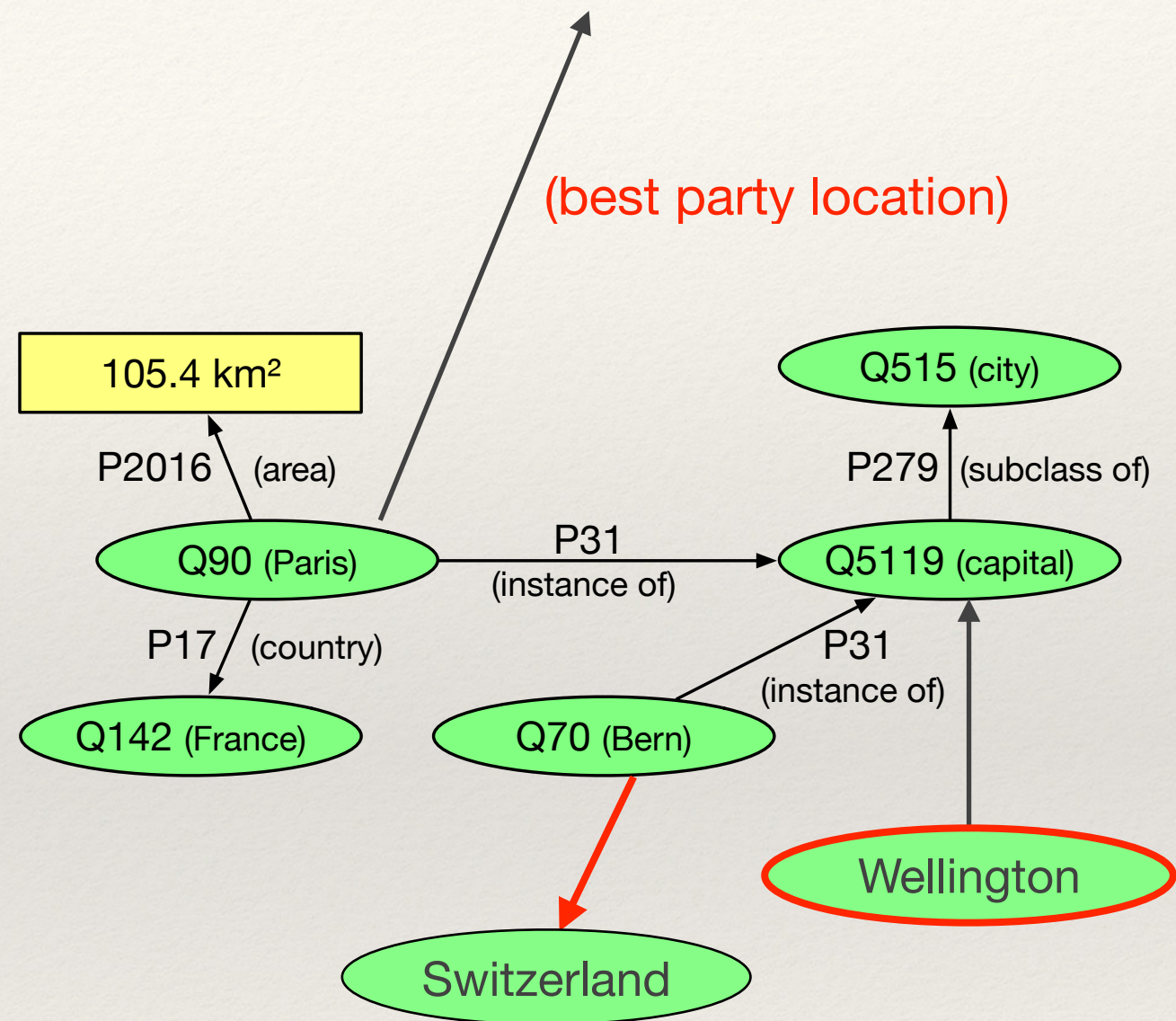- Class Completeness Estimators

- Evaluation / Application

# KG Completeness

Schema completeness

Property completeness

Interlinking completeness

*Class completeness*

(best party location)

105.4 km²

P2016  (area)

Q90 (Paris)

P17  (country)

Q142 (France)

P31
(instance of)

Q70 (Bern)

P31
(instance of)

Switzerland

Q515 (city)

P279  (subclass of)

Q5119 (capital)

Wellington

After: Zaveri, A., Rula, A., Maurino, A., Pietrobon, R., Lehmann, J., Auer, S.: Quality assessment for Linked Data: A survey. Semantic Web Journal

# Errors through incomplete Classes

Missing entities can lead to **wrong conclusions:**

*"There are no volcanos in New Zealand, so no need for an early warning system."*

Missing entities can **bias statistics**:

*"There are more Skyscrapers in Auckland, compared to NY, so Auckland is bigger."*

# The Question

How can we know if we have all real world entities of a class **C** in our Knowledge Base?

How many **Volcanos** are there?
How many **Hospitals** are there?
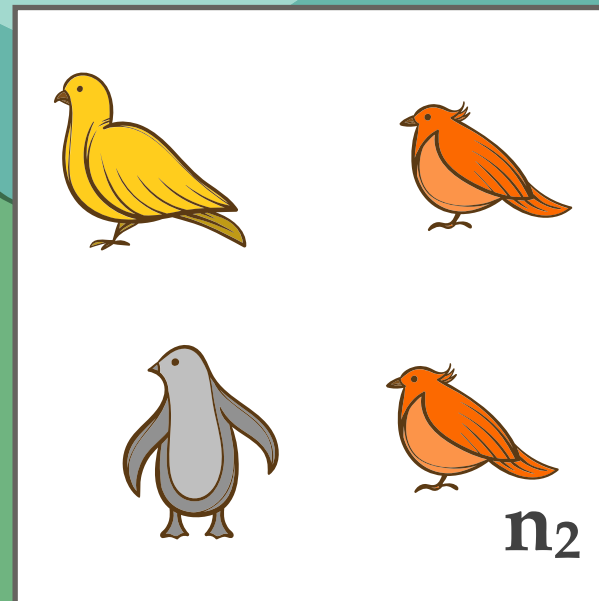How many **........** are there?

$$I_C = \{I_1, ..., I_N\}$$
$$N = |I_C|$$

How many **I** has $I_C$ ?

How many **Mountains** are there?

IMountains

# Species Richness Estimators

# Species Richness Estimators



|  | $n_1$ | $n_2$ | $n_3$ | $n_4$ |
|---|---|---|---|---|
|  | I |  | I |  |
|  | I |  | I | I |
|  | I | I | I |  |
|  | I |  | I | I |
|  |  | I |  | I |
|  |  | I |  | I |
| $f_1$ | 4 | 5 | 2 | 0 |
| $f_2$ | 0 | 1 | 3 | 3 |
| $f_3$ | 0 | 0 | 2 | 3 |

# Collaborative Knowledge Graphs

# Wikidata Edits

Edits: 161'445'153

Classes: 54'698

Q90 (Paris) —P31 (instance of)→ Q5119 (capital)

Q70 (Bern) —P31 (instance of)→ Q5119 (capital)

E1: Paris

E2: Bern

E3: Bern

E4: Paris

E5: Bern

E6: Berlin

Timeline

# Class Completeness Estimators

# Class Completeness Estimators

**Jack1**      **Jackknife Estimators**

**N1-UNIF**    **Sample Coverage and the Good-Turing Estimator**

**SOR**        **Singleton Outliers Reduction**

**Chao92**     **Abundance-based Coverage Estimator**

# Sample Coverage
## and the Good-Turing Estimator

$$\hat{N}_{\text{N1-UNIF}} = \frac{D}{\hat{S}} = \frac{D}{1 - \frac{f_1}{n}}$$

**D**   Distinct Entities

$$S = \sum \mathbb{1}[X_i > 0]$$

**S**   *Sample Coverage*
**N**   True class Size
$p_i$   Probability to Observe
$X_i$   Frequency of Observation

$$\hat{S} = 1 - \frac{f_1}{n}$$
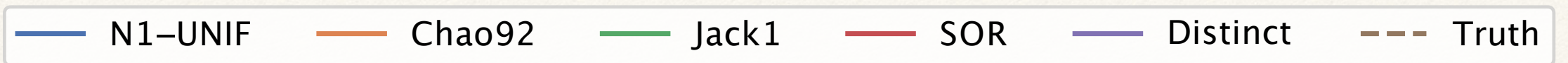
**Ŝ**   *Good-Turing Estimator*
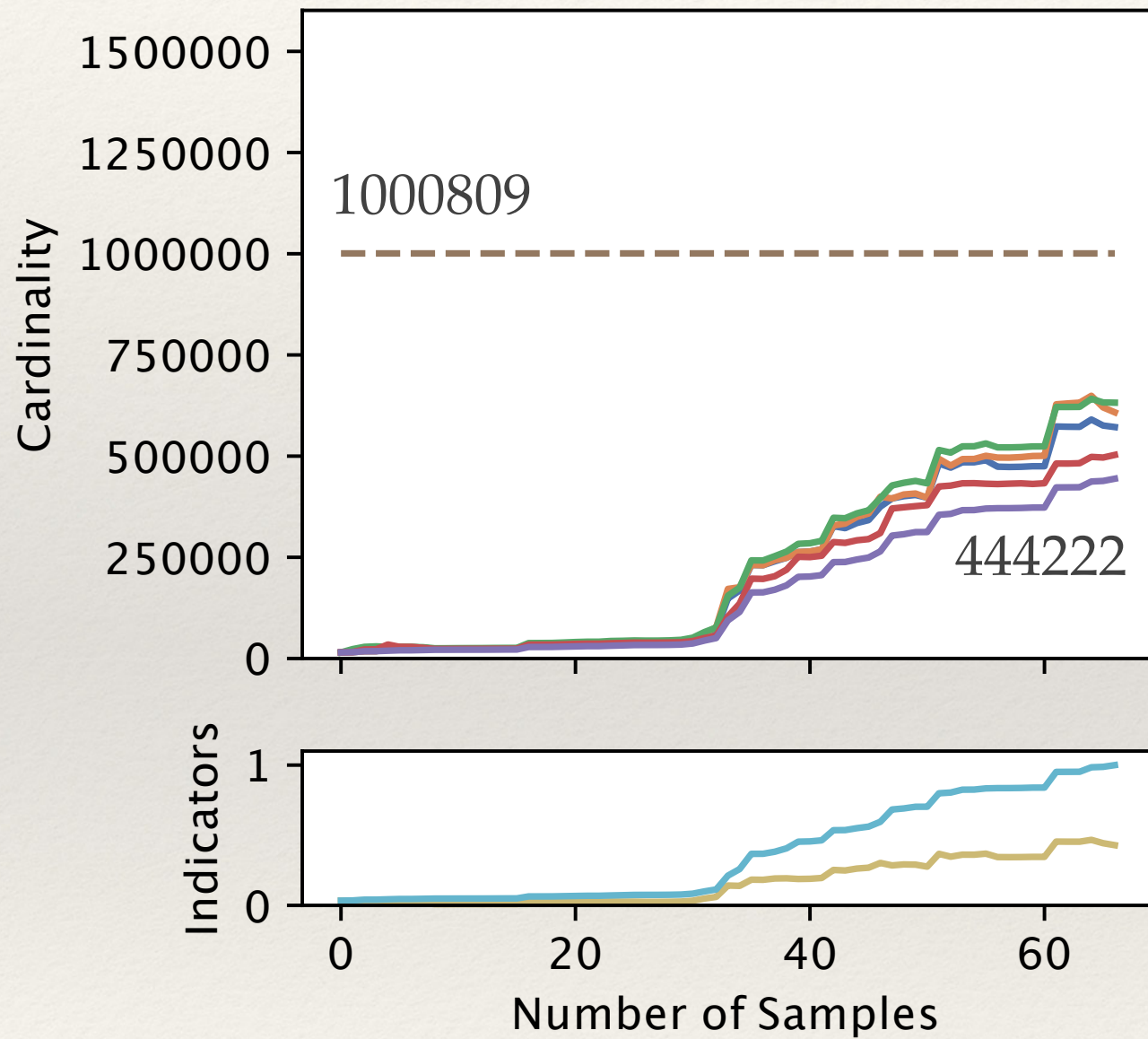$f_1$   Instances observed once
$n$   Number of Observations
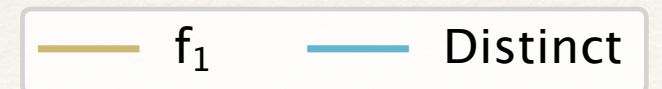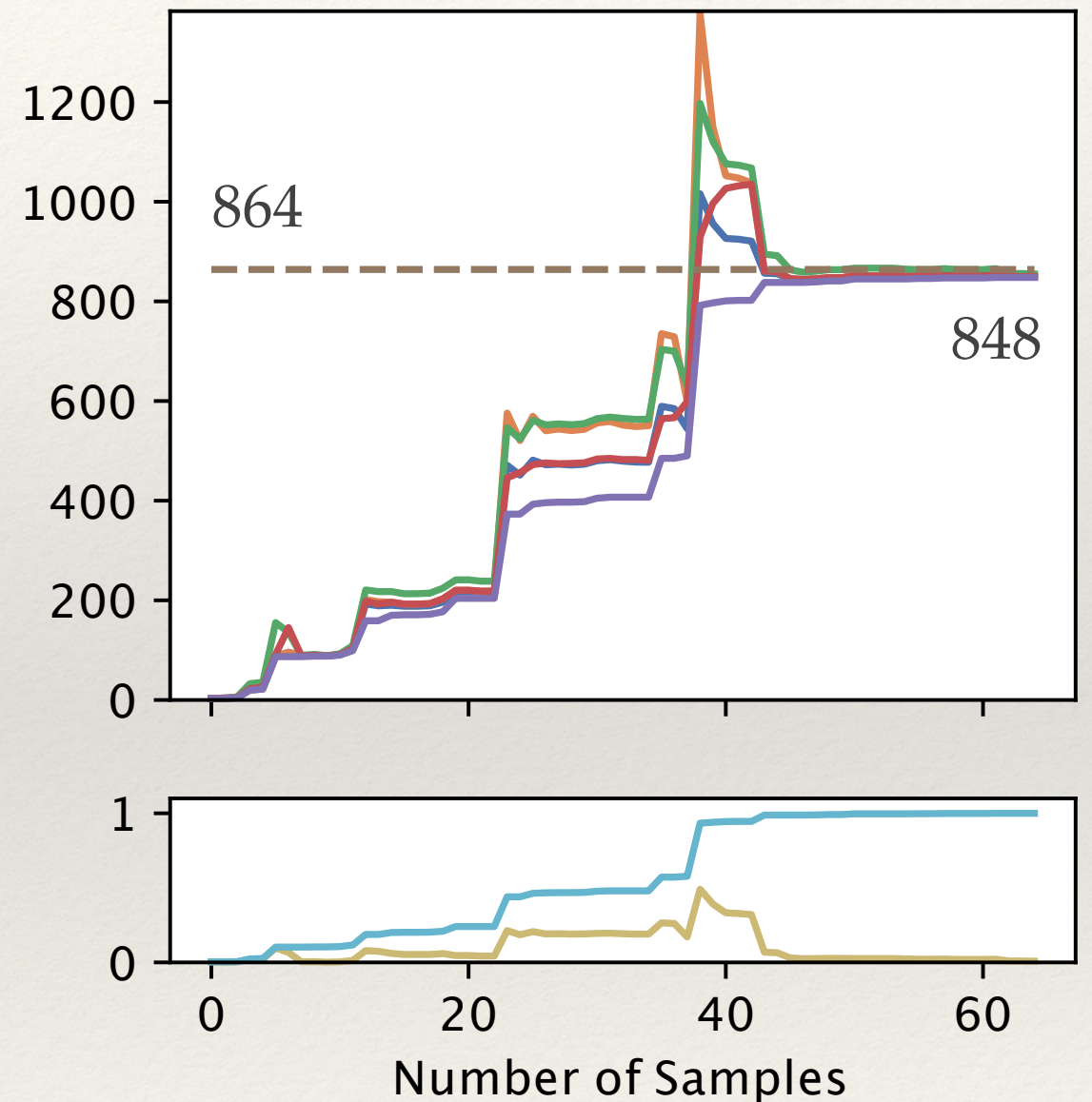
16

# Evaluation

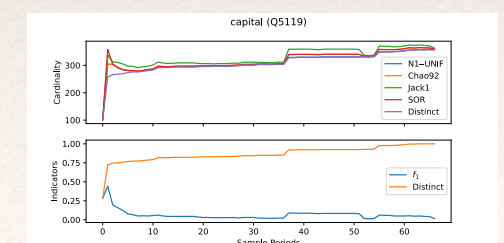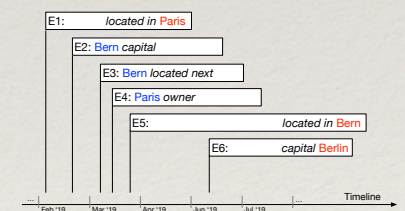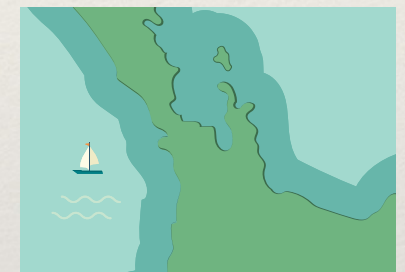# Evaluation

# Application
## Convergence Metric

$$\rho = \frac{\sum_{i=k-w}^{k} \frac{|\hat{N}_i - D_i|}{D_i}}{w}$$

$\hat{N}_i$    Entities Estimate per Period
$D_i$    Distinct Entities per Period
$w$    Window

| SOR $\rho < 0.001$ Distinct | | | SOR $\rho > 0.1$ Distinct | | |
|---|---|---|---|---|---|
| municipality of Japan | 0.0000 | 739 | urban beach | 0.1759 | 683 |
| Philippine TV series | 0.0009 | 822 | hydroelectric power station | 0.2975 | 2,936 |
| Landgemeinde of Austria | 0.0000 | 1,116 | aircraft model | 0.1800 | 3,919 |
| district of China | 0.0009 | 975 | motorcycle manufacturer | 0.1758 | 690 |
| nuclear isomer | 0.0002 | 1,322 | local museum | 0.1760 | 1,150 |
| international border | 0.0000 | 529 | waterfall | 0.1942 | 5,322 |
| commune of France | 0.0001 | 34,937 | race track | 0.2783 | 946 |
| village of Burkina Faso | 0.0005 | 2,723 | film production company | 0.2107 | 2,179 |
| supernova | 0.0005 | 5,906 | red telephone box | 0.3469 | 2,716 |
| township of Indiana | 0.0002 | 999 | mountain range | 0.2390 | 21,390 |

# Wrap-Up

- The edit history of a KG can be used to inform statistical methods adapted from species estimators.



- We evaluated the effectiveness of statistical methods to estimate the class size on repeated sampling.



- With the convergence metric we are able to distinguish between complete and incomplete classes in a KG.



**https://cardinal.exascale.info/**