



Chair of
Databases and
Information Systems



GouDa – Generation of universal Data Sets

Valerie Restat, Gerrit Boerner, André Conrad, Uta Störl

University of Hagen

12.06.2022, DEEM-Workshop

Content

Introduction

GouDa

Properties

Error types

Usage

Conclusion and Future Work

Analysis of Data Preparation Pipelines

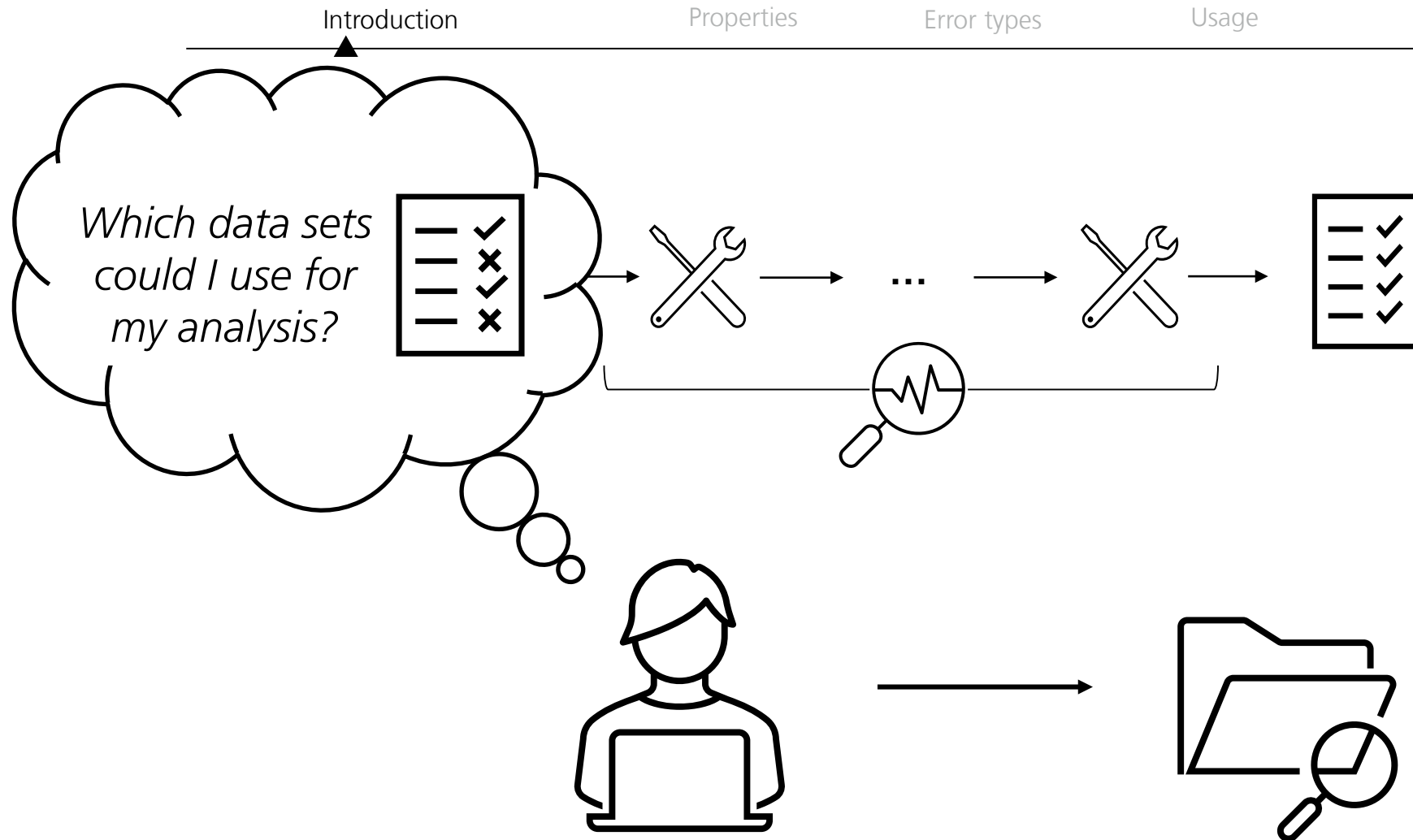
Introduction

Properties

Error types

Usage

Conclusion and Future Work



Data Sets

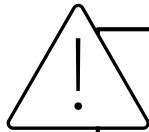
Introduction

Properties

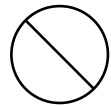
Error types

Usage

Conclusion and Future Work



Problem Statement 1



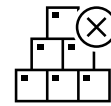
Not publicly available



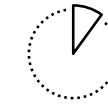
Missing ground truth



Problem Statement 2



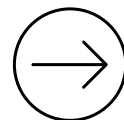
Error types



Error rates

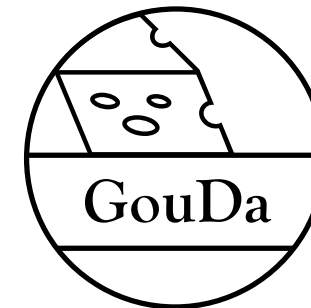


Size of data sets



GouDa – Test data generator

- Publicly available
- Ground truth provided
- Diverse error types
- Arbitrary error rates
- Scalable



Data Generators

Introduction

Properties

Error types

Usage

Conclusion and Future Work

- Different data generators exist: Software testing, SQL queries, ...
- Tools for generating data sets with defined errors: BART [Arocena et al. (2015)]
(Benchmarking Algorithms for data Repairing and Translation)



Errors are inserted into already existing, clean databases



Requires a relational database schema with unique tuple identifiers



Focusing on constraint-induced errors



No existing data needed

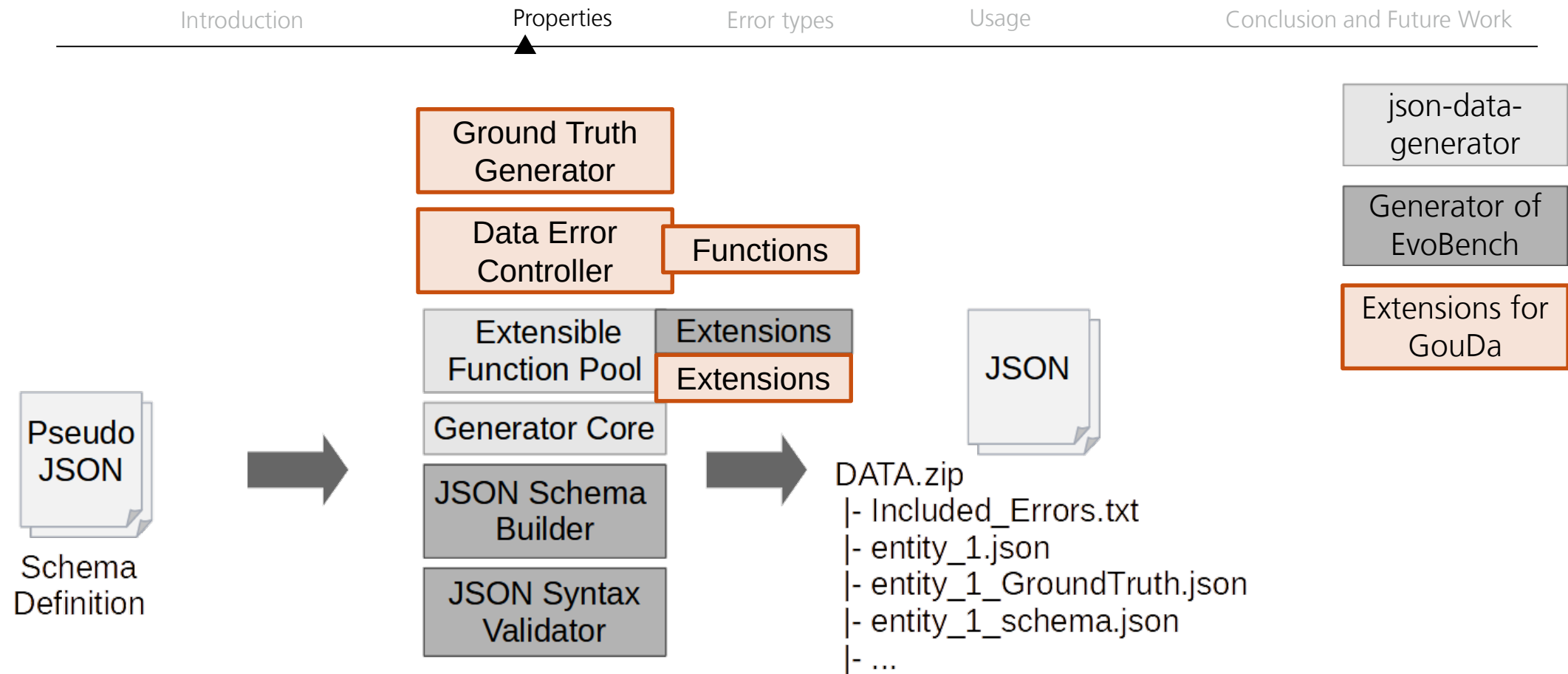


Uses JSON format, not only relational data but also nested data is supported



Supports variety of different error types

GouDa – Components



Conrad et al. (2021) *EvoBench: Benchmarking Schema Evolution in NoSQL.*,
 - based on json-data-generator: <https://github.com/vincentrussell/json-data-generator> (v1.12)

GouDa – Characteristics

Introduction

Properties

Error types

Usage

Conclusion and Future Work



Error rate



- *Freely configurable occurrence rate*
- *Errors are randomly distributed*

Reproducible

Scalable

Portable

Realistic

Ground truth

Error types

GouDa – Characteristics

Introduction

Properties

Error types

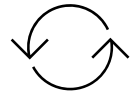
Usage

Conclusion and Future Work

Error rate



Reproducible



*Same type of error can be generated
on the same attribute of the same tuple*

Scalable

Portable

Realistic

Ground truth

Error types

GouDa – Characteristics

Introduction

Properties

Error types

Usage

Conclusion and Future Work

Error rate

Reproducible

Scalable

Portable

Realistic

Ground truth

Error types



Various options available

GouDa – Characteristics

Introduction

Properties

Error types

Usage

Conclusion and Future Work

Error rate

Reproducible

Scalable

Portable

Realistic

Ground truth

Error types

{ } *Data sets are generated in JSON format*

GouDa – Characteristics

Introduction

Properties

Error types

Usage

Conclusion and Future Work

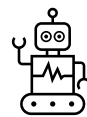
Error rate

Reproducible

Scalable

Portable

Realistic



Different lists with realistic names, words or addresses are included

Ground truth

Error types

GouDa – Characteristics

Introduction

Properties

Error types

Usage

Conclusion and Future Work

Error rate

Reproducible

Scalable

Portable

Realistic

Ground truth



Ground truth is provided

Error types

GouDa – Characteristics

Introduction

Properties

Error types

Usage

Conclusion and Future Work

Error rate

Reproducible

Scalable

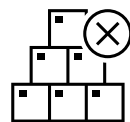
Portable

Realistic

Ground truth



Error types



Variety of different error types

GouDa – Error types

Introduction

Properties

Error types

Usage

Conclusion and Future Work

R	a ₁	a ₂	...	a _n
t ₁				
t ₂				
t ₃				
...				
t _m				

An Attribute Value
of a Single Tuple

R	a ₁	a ₂	...	a _n
t ₁				
t ₂				
t ₃				
...				
t _m				

The Values of a
Single Attribute

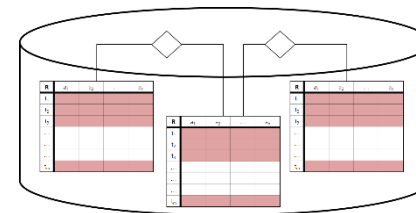
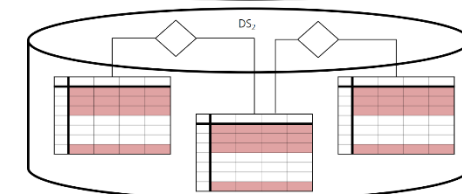
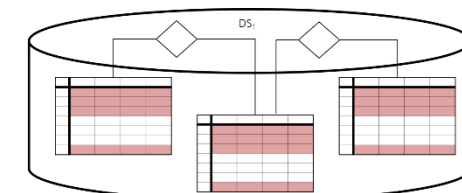
R	a ₁	a ₂	...	a _n
t ₁				
t ₂				
t ₃				
...				
t _m				

The Attribute Values
of a Single Tuple

R	a ₁	a ₂	...	a _n
t ₁				
t ₂				
t ₃				
...				
t _m				

The Attribute Values
of Several Tuple

Single Relation


 Relationships among
Multiple Relations


Multiple Data Sources

Single Source

Data Quality Problems

Oliveira et al. (2005)
*A taxonomy of data
quality problems*

GouDa – Error types

Introduction

Properties

Error types

Usage

Conclusion and Future Work

R	a ₁	a ₂	...	a _n
t ₁				
t ₂				
t ₃				
...				
...				
t _m				

An Attribute Value of a Single Tuple

- Missing value
- Syntax violation
- Interval violation
- Set violation
- Misspelled error
- Inadequate value to the attribute context
- Value items beyond the attribute context
- Meaningless Value
- Erroneous entry*

- Semi-empty tuple
- Inconsistency among attribute values
- Irrelevant observation*

R	a ₁	a ₂	...	a _n
t ₁				
t ₂				
t ₃				
...				
...				
t _m				

The Attribute Values of a Single Tuple

R	a ₁	a ₂	...	a _n
t ₁				
t ₂				
t ₃				
...				
...				
t _m				

The Values of a Single Attribute

- Uniqueness value violation
- Synonyms existence
- Outlier*
- Missing Attribute*

- Redundancy about an entity
- Inconsistency about an entity
- Bias*
- Noise*

R	a ₁	a ₂	...	a _n
t ₁				
t ₂				
t ₃				
...				
...				
t _m				

The Attribute Values of Several Tuple

GouDa – Usage

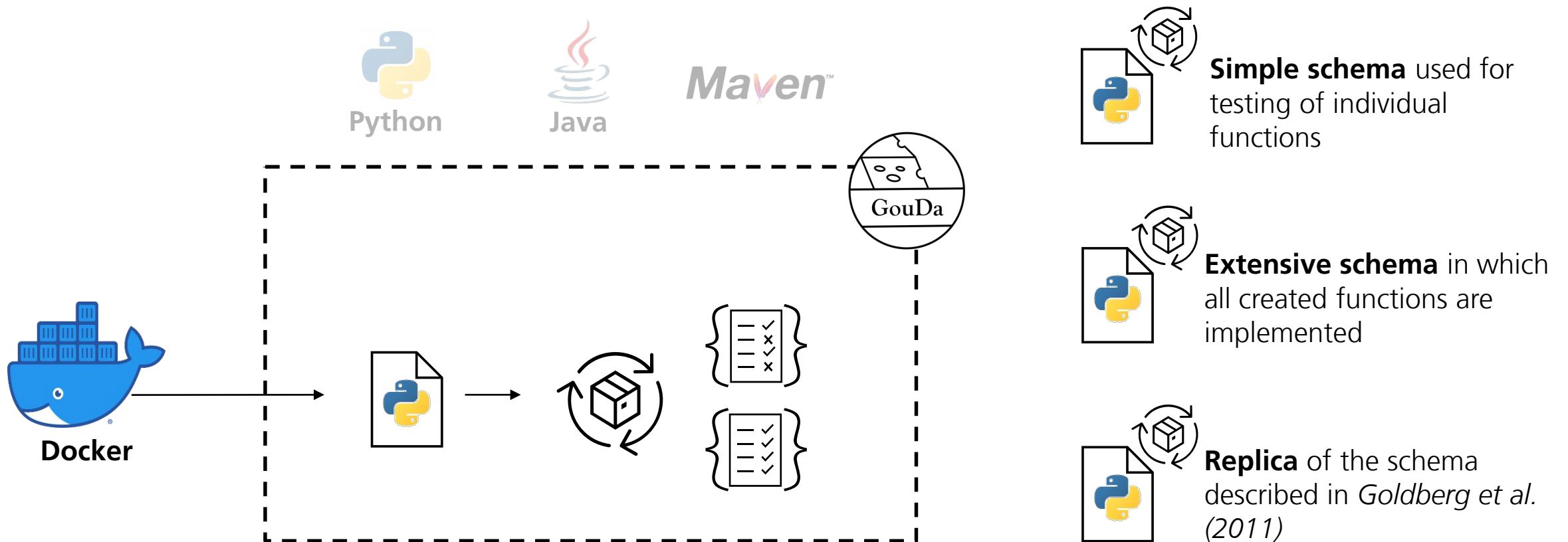
Introduction

Properties

Error types

Usage

Conclusion and Future Work



<https://zenodo.org/record/6610025>

GouDa – Usage

Introduction

Properties

Error types

Usage

Conclusion and Future Work

10000 data objects

```
[ '{repeat(10000)}', {
  "Index": {{relationship("Performance_Test-id-idx", "INCREMENT", "1")}},
  "City": "{city()}" Function for city names, no errors,
  "Date": "{error("SYNTAX", 5, "date", "dd.MM.yyyy", date("08-05-1945
00:00:00", "03-10-1990 00:00:00", "dd.MM.yyyy"))}" Date function, 5% syntax error,
  "Count": "{error("INTERVALLVIOLATION", 5, 0, 9999, integer(0, 9999))}" Integer range with 5% intervall violation
}]
```

Incremental Index

GouDa – Usage

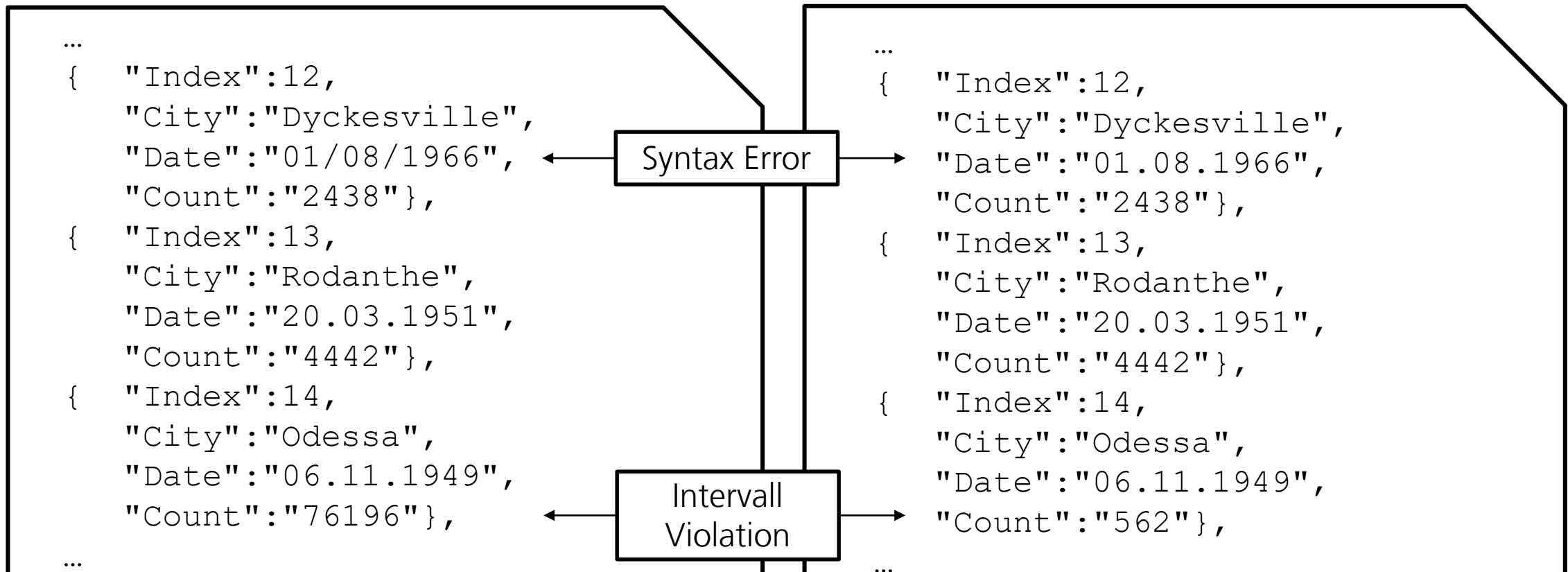
Introduction

Properties

Error types

Usage

Conclusion and Future Work

Data set with errors
Ground Truth


Conclusion and Future Work

Introduction

Properties

Error types

Usage

Conclusion and Future Work



We have presented **GouDa**, our test data generator



Diverse error types



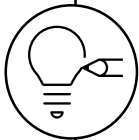
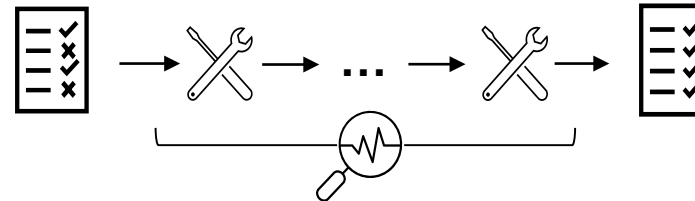
Arbitrary error rates



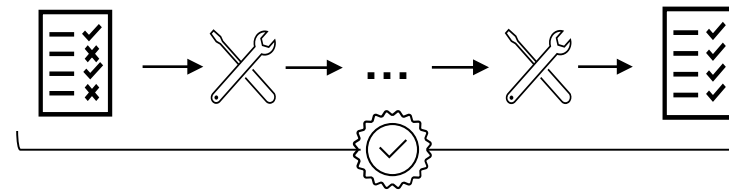
Provides ground truth



Which we currently use for: **Analysis** of data preparation tools and pipelines



And our next research topics: **Data quality** and **evaluation** of data preparation pipelines



GouDa – Generation of universal Data Sets

<https://zenodo.org/record/6610025>



Introduction

GouDa

Properties

Error types

Usage

Conclusion and Future Work

Valerie Restat

valerie.restat@fernuni-hagen.de

Gerrit Boerner

gerrit.boerner@studium.fernuni-hagen.de

André Conrad

andre.conrad@fernuni-hagen.de

Uta Störl

uta.stoerl@fernuni-hagen.de

Sources

- Patricia C. Arocena et al . 2015. Messing Up with BART: Error Generation for Evaluating Data-Cleaning Algorithms. Proc. VLDB Endow. (2015), 36–47. <https://doi.org/10.14778/2850578.285057>
- Paulo Oliveira et al. 2005. A taxonomy of data quality problems. In 2nd Int. Workshop on Data and Information Quality. 219–233.
- Saveli Goldberg, Andrzej Niemierko, and Alexander Turchin. 2008. Analysis of Data Errors in Clinical Research Databases. In AMIA 2008, American Medical Informatics Association Annual Symposium, Washington, DC, USA, November 8-12, 2008. AMIA. <https://knowledge.amia.org/amia-55142-a2008a-1.625176/t-001-1.626020/f-001-1.626021/a-049-1.626417/a-050-1.6264>