# DATA PREPARATION & CLEANING

## Chapter 2: Data Extraction

Victor Christen

## AGENDA

– General

Characteristics of Data Extraction

– Internal vs External

– ETL vs ELT

– Data Warehousing

– Extraction methods

   – API

   – Web Crawler & Scraping

   – Content specific extraction

   – Extraction from unstructured data

– Summary

# DATA EXTRACTION - GENERAL

"Data extraction is the process of **collecting** or **retrieving** disparate types of data from a **variety of sources,** many of which may be **poorly organized** or **completely unstructured**."[1]

- **Internal** vs **external** data source
- **ETL** vs **ELT**
- **Structured**
  - Standardized format, well defined access (query language)
- **semi/unstructured**
  - Websites, scanned and OCR-processed text documents, transcripted Audio Reports, etc.

Extraction methods depend on the data format and the type of data sources
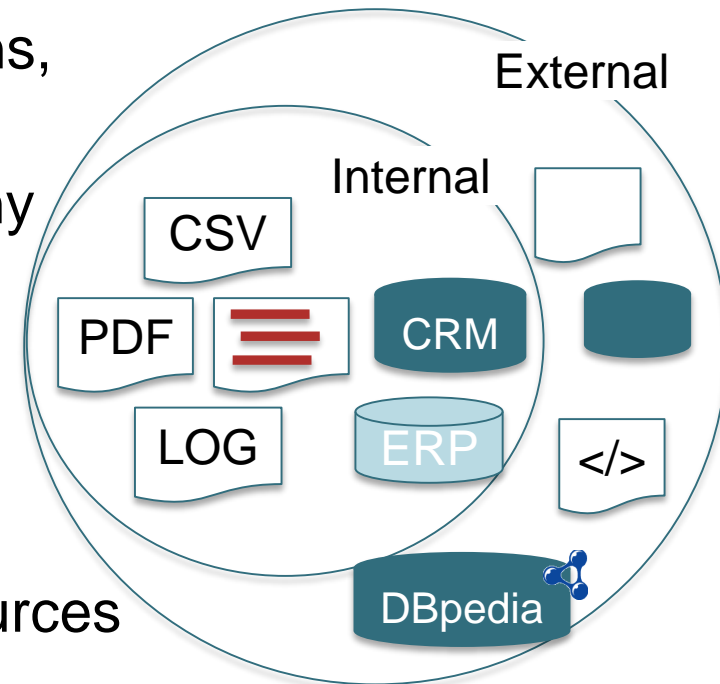
[1] https://www.talend.com/resources/data-extraction-defined/
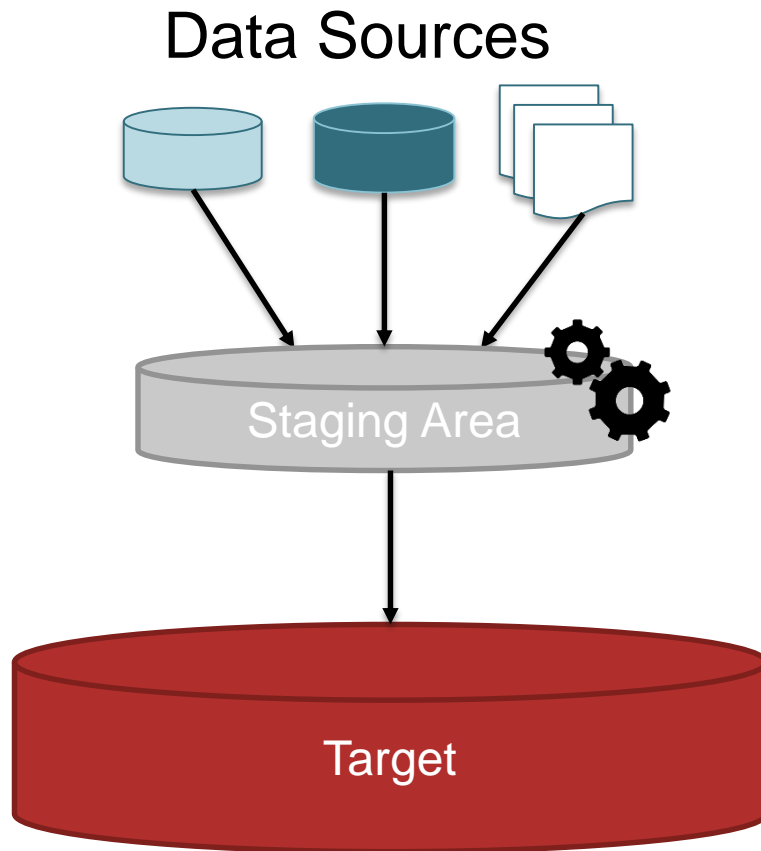
# EXTERNAL VS INTERNAL DATA SOURCES

## Internal

– Databases (CRM, ERP), documentations, Log-files, Emails

– View on internal processes of a company

– Completely autonomous w.r.t. the data source

– flexible access possibilities

## External

– Context information for internal data sources

– Websites, API, SPARQL-endpoints,…

– Access potentially limited by the functionality of an API, availability of resources

External

Internal

CSV

PDF

CRM

LOG

ERP

</>

DBpedia

# ETL VS ELT

## Data Sources



**Extraction**

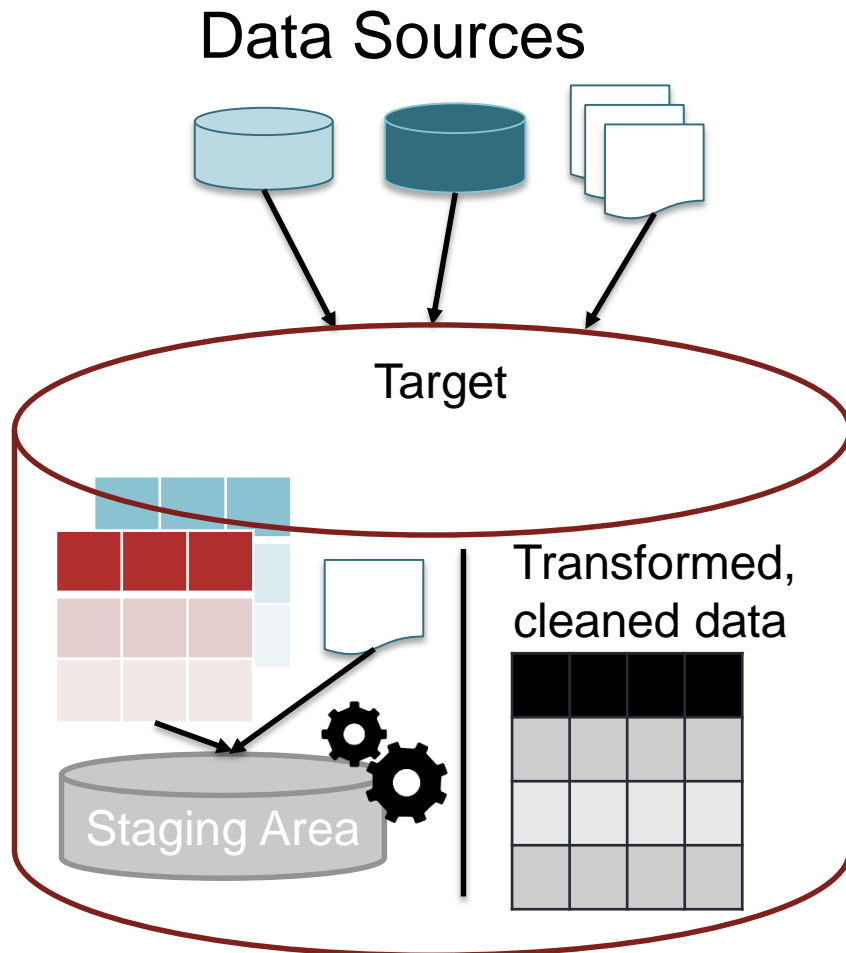– Extract data from data sources

**Transformation**

– Transform, clean and integrate data

**Loading**

– Load transformed, cleaned and integrated data in target database for analysis

# ETL VS ELT (2)



Data Sources

Target

Transformed, cleaned data

Staging Area

**Extraction**

– Extract data from data sources

**Loading**

– Load the extracted unprocessed data in the target data source, e.g., **Data Lake**

**Transformation**

– Transform and clean only the **relevant** data **before** analysis

## ETL VS ELT – PROS & CONS

## ETL

Target analysis tasks are known, e.g. product sales, human resources, production rate

+ Transformation before enables fast analysis

– Not scalable for large and complex transformation tasks

– Not flexible regarding changed analysis goals

## ELT

Direction of analysis is known but not in detail

+ Flexible regarding analysis tasks and data evolution

+ Store all data → flexible regarding new analysis tasks

o High performance using cloud infrastructure by using additional resources on demand

# ETL - DATA WAREHOUSING

– Specialized database to support company decisions

– Disconnected from operational systems

– Periodic ETL process

  – Support time dependent analysis

  – Initial loading after that mostly reading operations

–  Focus on specific target analysis such as product sales, citation numbers, …
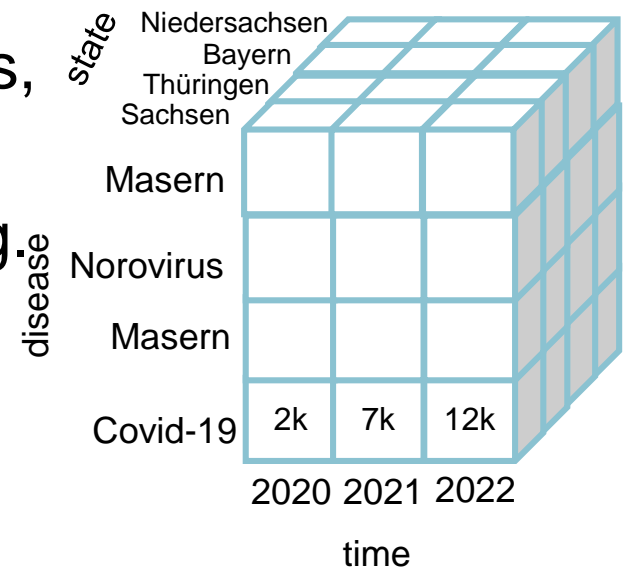
  – Described by different dimensions

## Reverse ETL

– Use Analysis result in operational databases

# DATAWAREHOUSE - COMPONENTS

- Data-Cube
  - Consisting of Dimensions and measures
- Measures are numerical values (counts, sales, etc.)
- Context information by dimensions (e.g. Time, place, name of disease)
- Dimensions represent coordinates for measures
- Hierarchical structure of dimensions possible
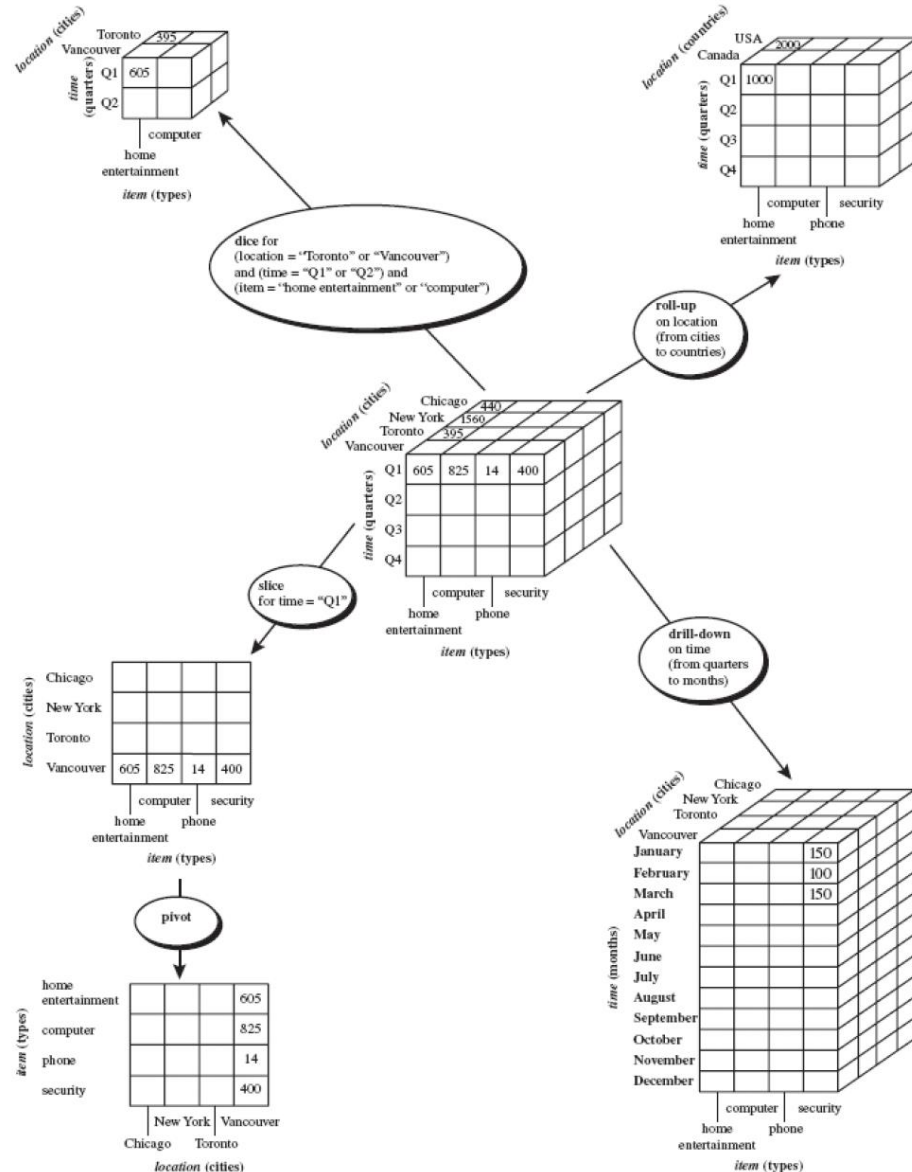
# DATAWAREHOUSING - OPERATIONS

Operators for analysis

- Roll-up (Aggregation of fine-granular hierarchies)
- Drill-down or Roll-down (detailed view)
- Slice and Dice (Restriction of one resp. multiple dimensions)
- Pivot (Switch orientation of dimensions)

Applications

- Generation of  statistics, reports, charts, etc.)
- Base to generate Data Mining models

Quelle: Han and Kamber, DM Book, 2 nd  Ed. (Copyright © 2006 Elsevier Inc.)

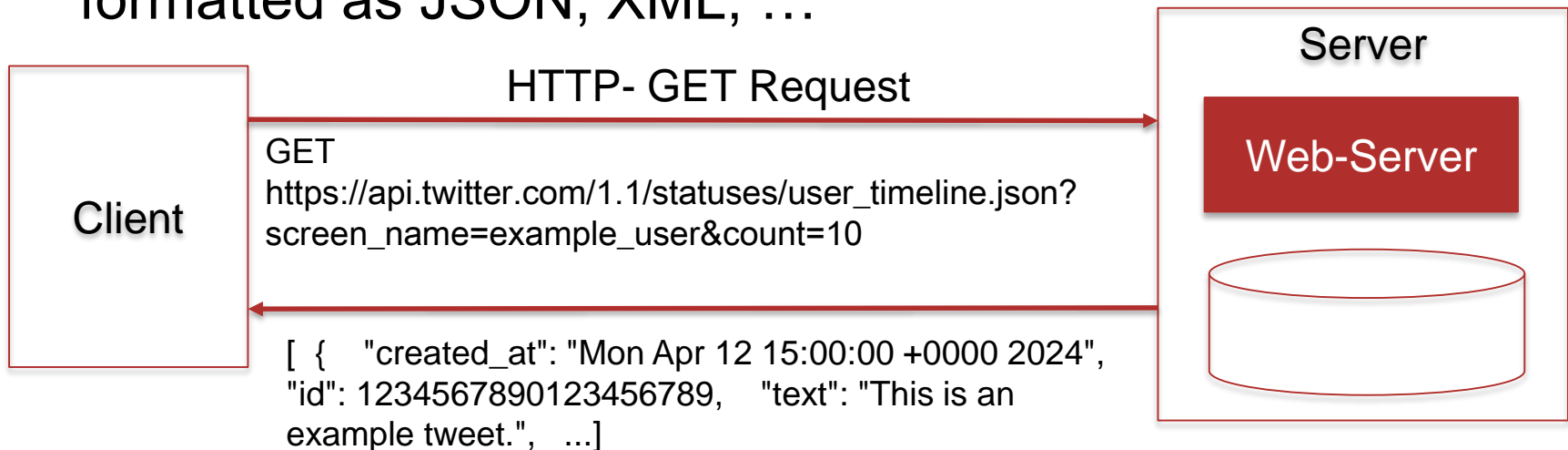# TYPES OF EXTRACTION METHODS

simple → complex

- Using a predefined Query Language
  - SQL, Cipher, SPARQL,
- Process relational data, e.g. CSV/TSV following a strict schema using program specific libraries such as Pandas, tool connectors

- Structure is defined by the data itself, e.g., XML/JSON
- Arbitrary set of attributes
- Potentially complex extraction methods required depending on the content of an attribute
- Query Language, API, Web Crawler/Scraping

- Data is completely unstructured, e.g. free formatted text, Images, Audio
- Requires Advanced methods such as Named Entity Recognition, OCR, speech-to-text
- Content specific libraries → tables

# REST - APPLICATION PROGRAMMING INTERFACE – REST API

– Website providers enable access to their resources by a REST-API

– Predefined set of access functionalities called by HTTP requests to retrieve the current state of resources

– Server responds by returning a semi-structured result formatted as JSON, XML, …

HTTP- GET Request

Server

Web-Server

Client

GET
https://api.twitter.com/1.1/statuses/user_timeline.json?
screen_name=example_user&count=10

[ {    "created_at": "Mon Apr 12 15:00:00 +0000 2024",
"id": 1234567890123456789,    "text": "This is an
example tweet.",   ...]

## API - GEOCODING

- Extraction of geo coordinates from address data for the visualization and analysis of location data
- Reverse determination of the corresponding address with regard to a geo coordinate

Applications

- Epidemiological research to determine disease spread clusters or analyze geographically related health problems

## GEOCODING - METHODS

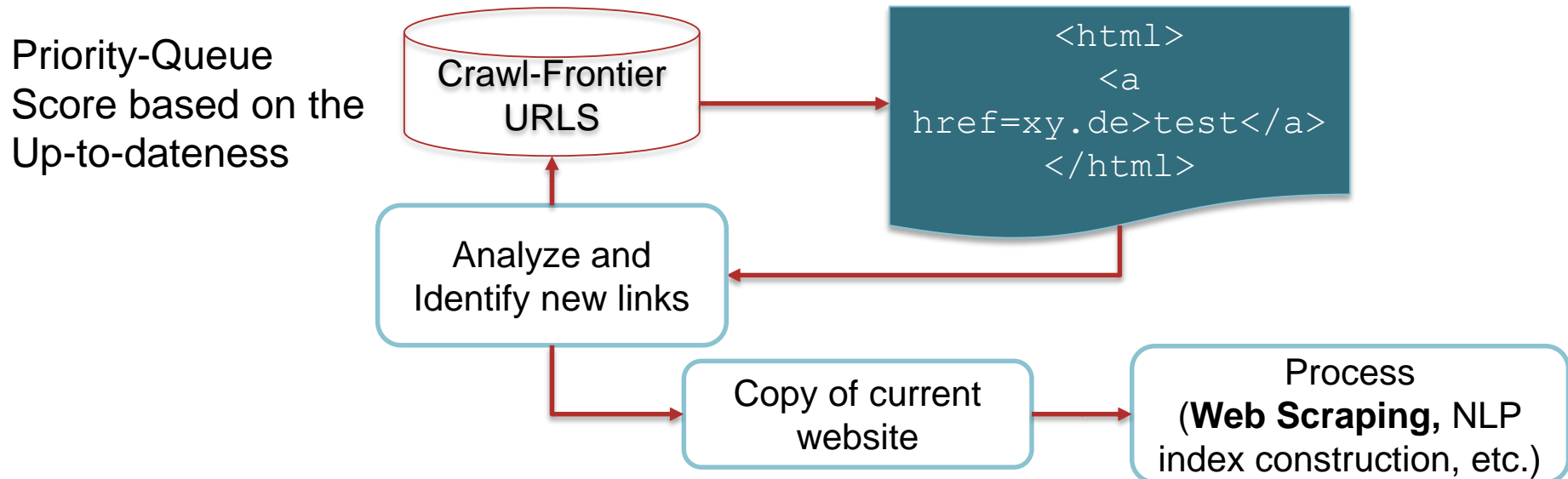– Use of a reference data set with corresponding mapping → often not available

API

– Standardized access to geo location data

– e.g. GeoPy, geocoder

– Exact match

  – 1:1 mapping address-coordinate

– Fuzzy Match

  – Specification of a region if there is no exact match with address

# DATA EXTRACTION FROM WEBSITES

## Web- Crawler

– Collect & Process **recursively** links and the content from websites based on the hyperlink structure

Priority-Queue
Score based on the
Up-to-dateness

Crawl-Frontier
URLS

```
<html>
  <a
href=xy.de>test</a>
</html>
```

Analyze and
Identify new links

Copy of current
website

Process
(**Web Scraping,** NLP
index construction, etc.)

– Periodical execution to notify changes of collected websites

# WEB SCRAPING

## Team

Home / Team

Current employees
Former employees

| Photo | Name | Type |
|---|---|---|
| | **PROF. DR. ERHARD RAHM** > | Employee |
| | ANDREA HESSE (SEKRETARIAT) > | Employee |

🔍 Search →

**Recent publications** 🔊

2024 / 4: **PHYSICS-INFORMED DEEP LEARNING TO QUANTIFY ANOMALIES FOR REAL-TIME FAULT MITIGATION IN 3D PRINTING** >

2024 / 3: **SERAPH: CONTINUOUS QUERIES ON PROPERTY GRAPH STREAMS** >

```
<div class="view-content">
<div class="w3-responsive">
<table class="w3-table-all cols-3 responsive-enabled">
    <thead>
        <tr>
            <th class="priority-medium views-field views-field-field-person-photo" id="view-field-person-photo-table-column"
            <th id="view-title-table-column" class="views-field views-field-title" scope="col"><a href="?field_person_type_v
            <th id="view-field-person-type-table-column" class="views-field views-field-field-person-type" scope="col"><a h
        </tr>
    </thead>
    <tbody>
        <tr>
            <td class="priority-medium views-field views-field-field-person-photo" headers="view-field-person-photo-table-co
                </a>
            </td>
            <td class="readon views-field views-field-title" headers="view-title-table-column"> <a href="/person/rahm" href
            <td headers="view-field-person-type-table-column" class="views-field views-field-field-person-type"> Employee
        </tr>
        <tr>
            <td class="priority-medium views-field views-field-field-person-photo" headers="view-field-person-photo-table-co
                </a>
            </td>
            <td class="readon views-field views-field-title" headers="view-title-table-column"> <a href="/person/andrea_hess
            <td headers="view-field-person-type-table-column" class="views-field views-field-field-person-type"> Employee
        </tr>
```

# WEB SCRAPING - LIBRARIES

− Extraction of specified elements via libraries, e.g. Scrapemark[1], Scrapy[2], BeautifulSoup[3] and the definition of website-specific element patterns

```python
from bs4 import BeautifulSoup
import requests

webFile = requests.get("https://dbs.uni-leipzig.de/de/person")
soup = BeautifulSoup(webFile.content, 'html.parser')
first_table = soup.select_one("div.w3-responsive table.w3-table-
all.cols-3.responsive-enabled")
list = soup.select(
"tr:has(> td.readon.views-field.views-field-title) a")
```

− `soup.tag` returns the selected value with tags, without by using `.string`

− `element.find_all(element_name)` returns a list of specified elements regarding the selected ancestor element

[1] https://github.com/arshaw/scrapemark
[2] https://scrapy.org/
[3] https://www.crummy.com/software/BeautifulSoup/bs4/doc/

# WEB SCRAPING – CHALLENGES & RESTRICTIONS

- Evolution of websites → adaption of patterns
- Classical Scraping limited to static content, JavaScript content require evaluation
- Anti-Scraping mechanisms have to be considered and must not be circumvented
  - Access limitation (number of requests for a certain time interval)
  - Captcha
  - IP-Blocker

Legal restrictions

- Guarantee of copyrights depending on the use case → research vs commercialised application

UNIVERSITÄT LEIPZIG

# WEB SCRAPING - TOOLS

**User-Friendly Interface**

– Specify relevant elements through a graphical user interface

**Proxy Support**

– Avoid IP bans by proxies and rotation mechanism

**JavaScript Rendering**

Extraction of JavaScript generated content requires rendering functionalities

**Scheduled Scraping**

– Definition of crawling/scraping intervals

**Compliance and Ethical Considerations**

– Ensure compliance with website terms of service, legal regulations, ethical considerations, e.g., respecting robots.txt directives, avoiding excessive requests

# EXTRACTION FROM SPECIFIC CONTENT TYPE

## Tabular data

– Programming language specific libraries ,e.g., `tabula-java, tabula-py`

  – Automatic extraction of table data in a PDF-document

  – specification of pages, area, output options



PDF

Programming language specific table representation

```
import tabula
pdf_path =
"https://github.com/chezou/tabula-
py/raw/master/tests/resources/data.pdf"

dfs = tabula.read_pdf(pdf_path,
stream=True)
# read_pdf returns list of DataFrames
```

# EXTRACTION FROM UNSTRUCTURED DATA

Manual

- Handcrafted rules and background knowledge, e.g. Extracting author information from publications

- High quality of extracted data being correct

- Potential high number of missing data

  - Depending on the heterogeneity

- Very sensitive regarding changing content structure → High adaption effort

– Combine multiple extraction methods tailored to the certain task



| Named Entity Recognition (NER) | Entity Linking |
|---|---|
| • Identification of text boundaries of an named entity <br> • Utilizing domain-specific NER-model trained on a corpus from the same domain | • Link identified mentions to standardized entities of a common knowledge base such as ontology |

# AUTOMATIC EXTRACTION - EXAMPLE

– Processing medical documents such as electronic health record

– Enhance comparability by NER & EL

Standardized Entities



| Unfied Medical Language System (UMLS) |
|---|
| C0701055 |
| C0162723 |
| C0360120 |
| C0594492 |

# SUMMARY

- Internal vs External data sources
- ETL vs ELT
    - ETL: for small projects, complex transformation tasks potentially decrease the efficiency depending on the amount of data and heterogeneity, e.g., Data warehousing
    - ELT: transformation on demand considering a subset, High scalability using cloud environments, e.g., Data Lakes
        - increase resources on demand
- Data Warehousing
- Various extraction methods depending on the data source and format
    - Querying, processing relational data
    - API data extraction by predefined functions, e.g., Geocoding
    - Web crawler/web scraping
    - Content depending extraction methods, tabular data
    - Unstructured data requires complex and sophisticated methods such NER and entity linking