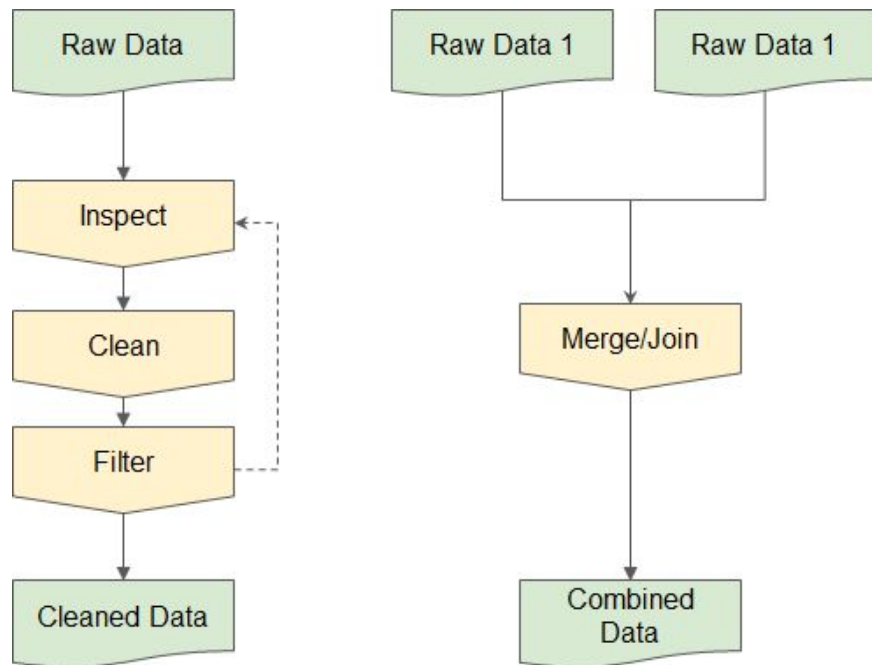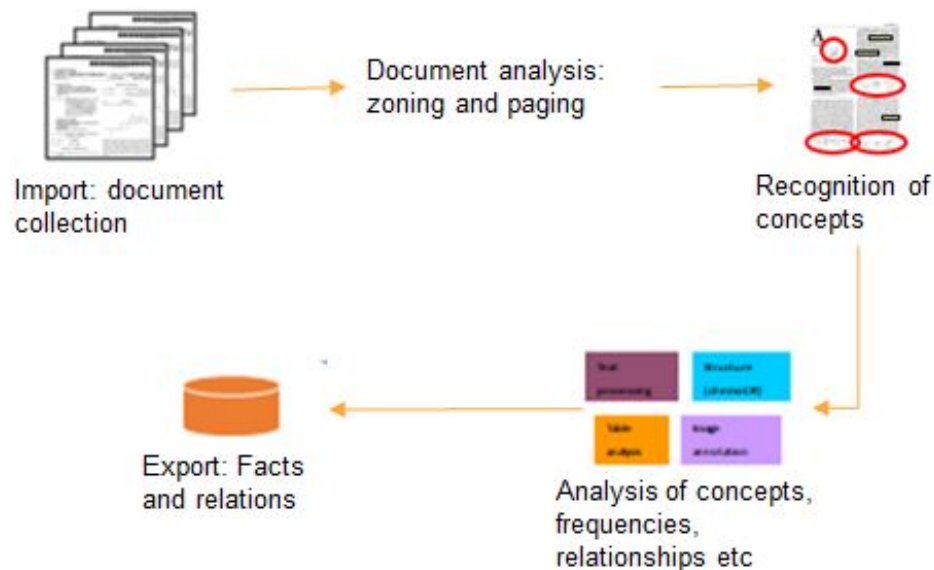# OpenRiskNet

# DataCure
## Data curation and creation of pre-reasoned datasets and searching

Noffisat Oki, Tim Dudgeon, Marc Jacobs, Danyel Jennen, Thomas Exner
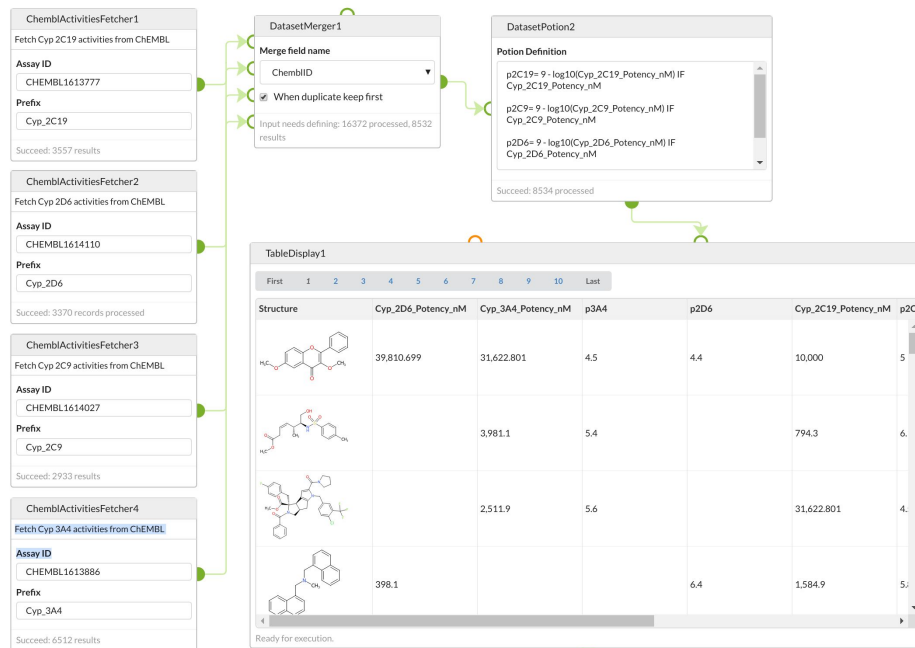
# Case Study objective



Data curation and merging

Text mining

# CypP450 data curation with Squonk

- Merge multiple datasets from ChEMBL into single set
- Uses ChEMBL identifiers to identify common structures
- Generates a dataset that can be used for machine learning
- See on [GitHub](GitHub)

# Data merging via data APIs

# OpenAPI + JSON-LD

```
"/datasets/749ea637-1262-46cc-b7cd-a8c2c6004245/data": {
    "get": {
        "operationId": "dataset",
        "x-orn-@id": "https://datasets/802a506e-2e08-4baf-b4f0-c3b89cbddd4d/data",
        "x-orn-@type": "x-orn:Service",

        "x-orn:expects":
            {
            },
        "x-orn:returns":
            {
                "x-orn-@id": "x-orn:ToxCastRawDataset"
            },
```

```
"x-orn-@context": {
    "x-orn-@id": "@id",
    "x-orn-@type": "@type",
    "@vocab": "http://openrisknet.org/schema#",
    "x-orn": "http://openrisknet.org/schema#",

    "unit":            "http://purl.obolibrary.org/obo/UO_0000000",
    "id":              "http://edamontology.org/data_0842",
    "SampleID":        "http://edamontology.org/data_3273",
    "Substance":       "http://ncicb.nci.nih.gov/xml/owl/EVS/Thesaurus.owl#C45306",
    "Compound":        "http://ncicb.nci.nih.gov/xml/owl/EVS/Thesaurus.owl#C43366",
    "PubChemCID":      "http://semanticscience.org/resource/CHEMINF_000140",
    "CompoundName":    "http://semanticscience.org/resource/CHEMINF_000043",
    "InChIKey":        "http://semanticscience.org/resource/CHEMINF_000059",
    "SMILES":          "http://semanticscience.org/resource/CHEMINF_000018",
    "CAS":             "http://semanticscience.org/resource/CHEMINF_000446",
    "Concentration":   "http://ncicb.nci.nih.gov/xml/owl/EVS/Thesaurus.owl#C41185",
    "Activity":        "http://www.bioassayontology.org/bao#BAO_0080024"
}
```

Subject or object

Predicate

# Finding datasets

## Selection of data service providing the XTT Cytotoxicity assay via API

```
[2]: sparql = SPARQLWrapper("http://orn-registry-openrisknet-registry.prod.openrisknet.org/api/sparql")
     serviceQuery = '''
     PREFIX orn: <http://openrisknet.org/schema#>
     SELECT *
     WHERE {
     ?tool orn:info ?info .
     ?info orn:title ?title .
     ?s1 <http://www.w3.org/1999/02/22-rdf-syntax-ns#type> <http://openrisknet.org/schema#EdelweissDataSet>
     }
     '''

     sparql.setQuery(serviceQuery)
     sparql.setReturnFormat(JSON)
     sparql_result = sparql.query().convert()
```

```
[3]: EdelweissServices = pd.DataFrame()

     for result in sparql_result:
         queryresults = pd.DataFrame(result['Result'][1]['ResultValues'], columns=result['Result'][1]['Variables'])
         EdelweissServices = pd.concat([EdelweissServices, queryresults])

     EdelweissServices[['title']]
```

```
[3]:                                               title
     0   ATG_XTT_Cytotoxicity_up_raw with OpenAPI schem...
```

```
[4]: ORNservice = "ATG_XTT_Cytotoxicity_up_raw with OpenAPI schemas and predicate annotation"
```

# Text mining

API

K8s hosted:
- https://api.scaiview.com/swagger-ui.html
- https://sam1.api.scaiview.com/swagger-ui.html

Workflows

jupyter https://biohub.scai.fraunhofer.de



**Biomedical Knowledge Discovery Pipeline**

**Public data**
· PubMed
· Articles
· Experimental data

BEL IEF
PRO MINER

Document storage

SCAI VIEW
BEL COMMONS

New ideas
Label data
FAIR data
Export data
Share data

Data Sources ⟩ Ingest ⟩ Storage ⟩ Discovery ⟩ Reporting & Visualization

Upload your own data

Extract
Transform
Load

TEM OWL
AETIONOMY KNOWLEDGEBASE

Search
Access
Analyse
Add

STUDY VIEWER
INTEGRATIVE VIEWER

https://biohub.scai.fraunhofer.de

Files  Running  Clusters

Select items to perform actions on them.

Upload  New

0 ▾  / home / marc

Name ↓  Last Modified  File si...

⌂ ..  seconds ago

▢ scaiview.ipynb  Running 5 hours ago  142

▢ Untitled.ipynb  Running 25 days ago  1.63

▢ untitled.txt

**swagger**

**SCAIView API** 2.0

[ Base URL: api.scaiview.com/ ]

https://api.scaiview.com/v2/api-docs?group=full-scaiview-api

SCAIView is a repository for biomedical documents that aims to provide semantic biomedical annota...

Terms of service

Apache License Version 2.0

**acl-controller** ACL Controller

**corpus-controller** Corpus Controller

**document-controller** Document Controller

**entity-statistics-controller** Entity Statistics Controller

**search-controller** Search Controller

POST /api/v5/corpora/{corpusID}/search  some nice desc.

POST /api/v5/corpora/{corpusID}/search/documents  some nice desc.

POST /api/v5/corpora/{corpusID}/search/export  some nice desc.

POST /api/v5/corpora/{corpusID}/search/identifier  some nice desc.

**solr-controller** Solr Controller

**Models**

Home  ✕  scaiview  ✕

← → ⟳  localhost:8888/notebooks/scaiview.ipynb

**jupyter** scaiview Last Checkpoint: Last Friday at 14:31 (autosaved)  Logout

File  Edit  View  Insert  Cell  Kernel  Widgets  Help  Trusted | Python 3 ○

Markdown

```
In [157]: def semanticQuery(concepts =[ human , cancer ], recall =true ):
          "create a fulltext query"

          if recall==True:
              operator = "OR"
          else:
              operator = "AND"

          query = {
              "operator": operator,
              "searchedConcepts": concepts
          }

          return query

In [155]: def fetchDocumentsByQuery( query, corpusID, limit='10' ):
          "fetch documents by query return json array"

          try:
              url = scaiview_uri+'corpora/'+corpusID+'/search/documents?size='+limit
              payload = json.dumps(query)
              headers = {'Accept': '*/*', 'Authorization': token, 'Content-Type': 'application/json'}

              r = requests.post(url, data=payload, headers=headers)
              documents = r.json()

              print ('got: ' documents['totalElements'])
              return documents

          except:
              print(payload)
              print(r.text)
              return ""

In [182]: def getTitleOfDocument(documents, i=0):
          "return the title of a document in a result list"

          return documents['content'][i]['documentElement']['metaElement']['bibliographic']['title']['titleText']['text']

In [183]: getTitleOfDocument(fetchDocumentsByQuery(fulltextQuery(['cerebrospinal fluid','brain'], False),corpusID,'1'))
          19458
Out[183]: '[Progressive multifocal leukoencephalopathy developing subsequent to cord blood transplantation in a patient with
          severe aplastic anemia].'

In [184]: getTitleOfDocument(fetchDocumentsByQuery(semanticQuery(['mesh:D002555','mesh:D001921'], False),corpusID,'1'))
          1248
Out[184]: 'The neuropathology of chronic traumatic encephalopathy.'

In [ ]: fetchDocumentAndRenderHTML('sha512Hex:24c1378
          7e8b9b83128bc930c3b27f3f5f2ae4e44677222833466a2e705cd0e1d44f4125bf652102c42bf0a667bd7e956
```

# OpenRiskNet example workflow

Task:

- Identify the concept of acetaminophen (definition, identifiers, synonyms)
- Find all relevant documents in the context of acetaminophen and carcinogenity
- What are the most relevant statements

Technology:

- Semantic index of PubMed/PMC (> 20 terminologies)
- Solr index + OLS index + UIMA pipeline