



ScholarSuccess

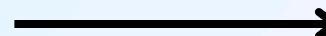
An End-to-End Pipeline for Conference Acceptance rate Prediction

A Project for 2110403 Introduction to Data Science and Data Engineering

CΣDT



WHY CONFERENCE PROCEEDING ?

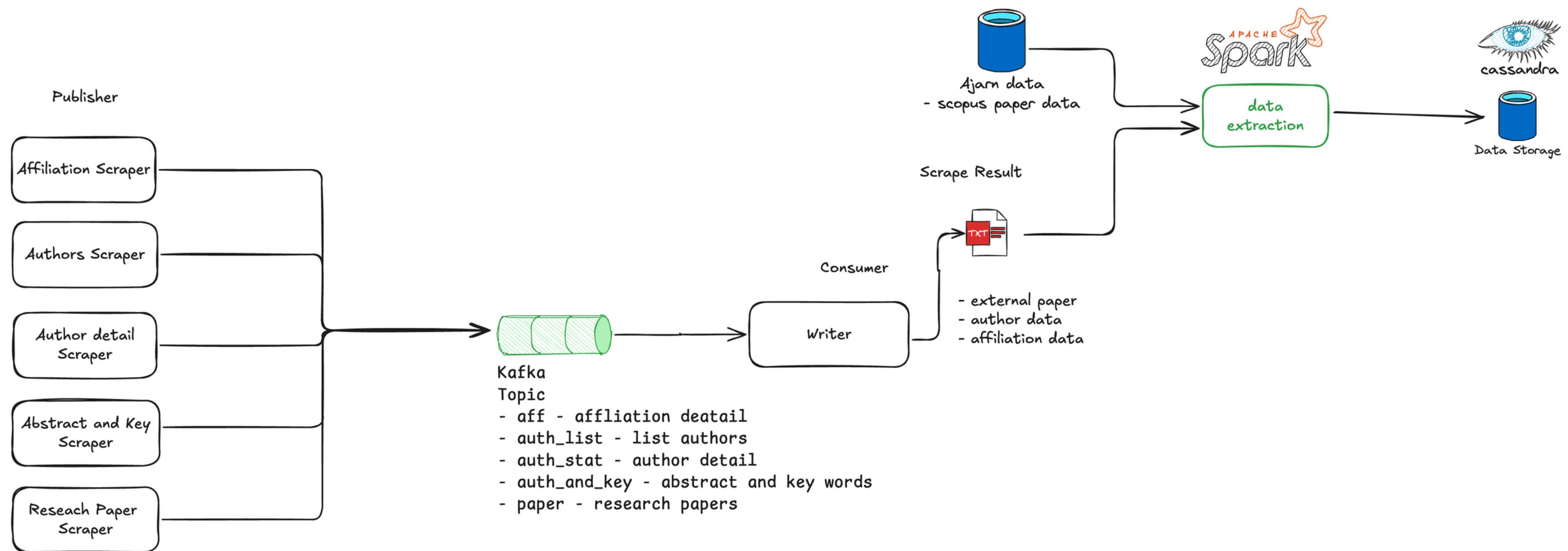


- Search for an important and efficient way to make others aware of a researcher's work and have it visible faster when an academic conducts a literature review
- help make researchers aware of new insights and industry innovations, while citing research underway on a given subject



DATA ENGINEERING

Data Engineering



DATA COLLECTION



- Data Provided (2018-2023)

```
{"abstracts-retrieval-response": {  
    "item": { ...  
    },  
    "affiliation": [ ...  
    ],  
    "coredata": { ...  
    },  
    "idxterms": {"mainterm": [ ...  
    ]},  
    "language": {"@xml:lang": "eng"},  
    "authkeywords": null,  
    "subject-areas": {"subject-area": [ ...  
    ]},  
    "authors": {"author": [ ...  
    ]}  
}}
```

20216 records

- Data Scraped
 - Research Papers
 - Authors List of Research Papers
 - Author Details
 - Abstract, Reference, Keywords of Research Papers
 - Affiliation Details



Scopus®



Elsevier Developer Portal

DATA COLLECTION

- Scraping Strategy

```
import concurrent.futures

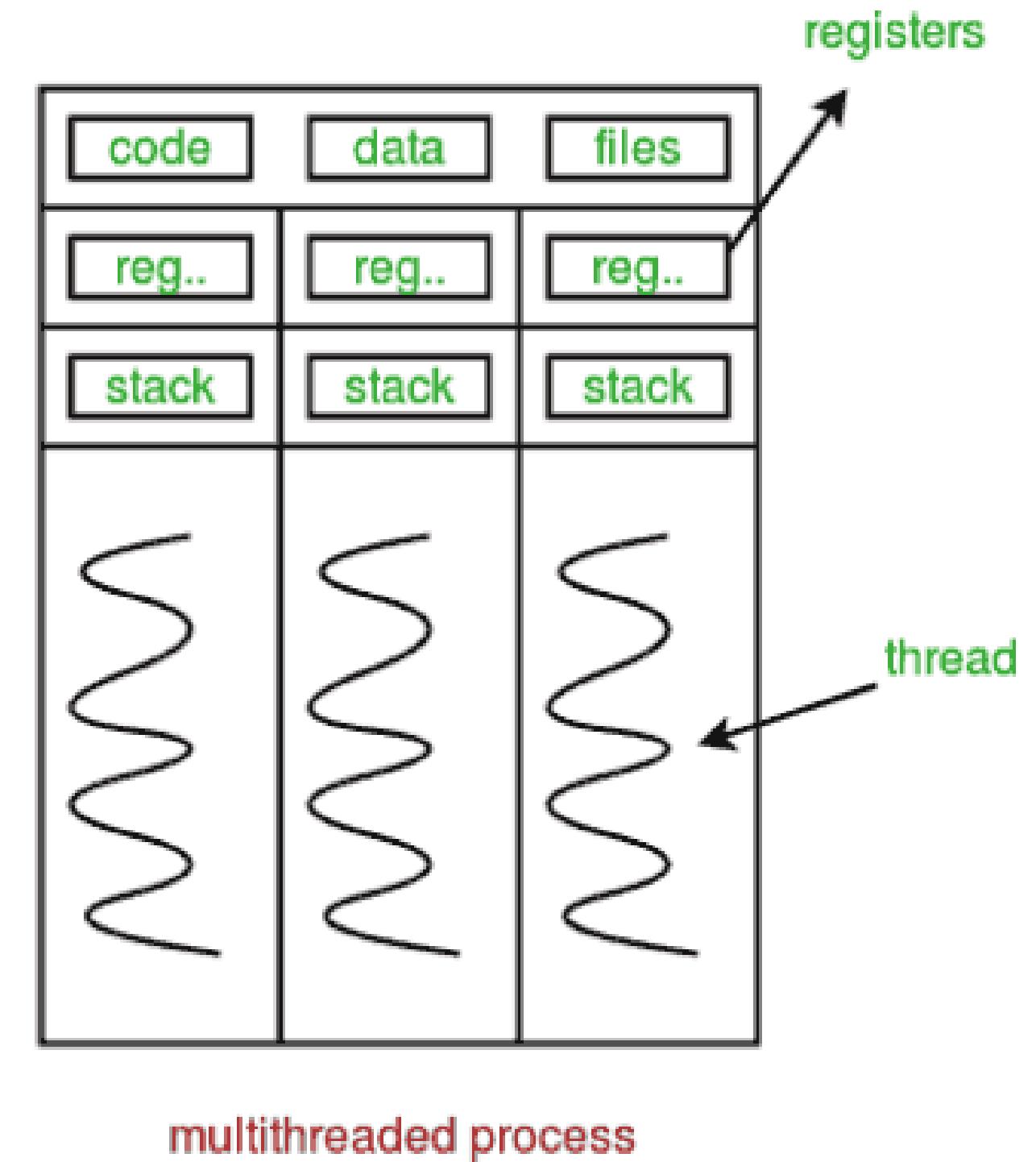
with concurrent.futures.ThreadPoolExecutor(max_workers=self.max_workers) as executor:
    author_results = list(filter(None, executor.map(self.scrape_author_data, authors)))
    :
    :
```

```
from concurrent.futures import ThreadPoolExecutor, as_completed

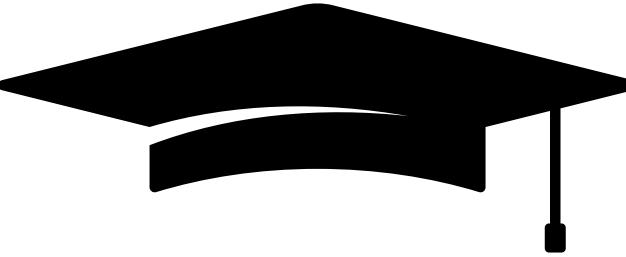
with ThreadPoolExecutor(max_workers=self.max_workers) as executor:
    progress_bar = tqdm.tqdm(
        total=len(self.data),
        desc="Fetching Abstracts",
        unit="paper"
    )

    futures = {
        executor.submit(self._process_paper, paper): paper
        for paper in self.data
    }

    for future in as_completed(futures):
        result = future.result()
        .
```



DATA COLLECTION



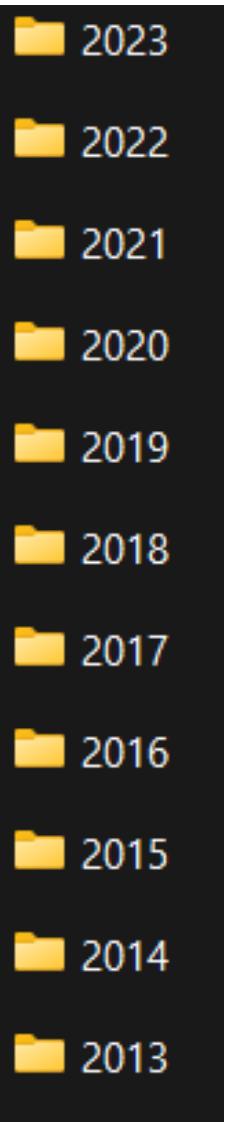
- Data Scraped
 - Research Papers
 - Authors List of Research Papers

42565 records

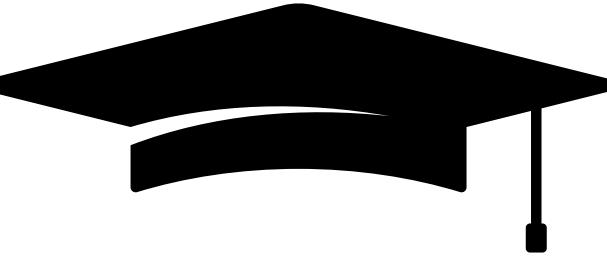
```
headers = {  
    'X-ELS-APIKey': self.api_key,  
    'Accept': 'application/json'  
}  
  
params = {  
    "query": queries[query_option],  
    "start": start_index,  
    "count": fetch_count,  
    "field": ...  
}  
  
"https://api.elsevier.com/content/search/scopus"  
  
response = requests.get(base_url, headers=headers, params=params)  
  
response = self.session.get(url, headers=headers, timeout=10)
```

Scopus API
➡ ➡ ➡

```
{  
    "@_fa": "true",  
    "prism:url": "https://api.elsevier.com/content/abstract/scopus_id/84891140835",  
    "dc:identifier": "SCOPUS_ID:84891140835",  
    "eid": "2-s2.0-84891140835",  
    "dc:title": "Vice chairman's message",  
    "dc:creator": "Gulati S.",  
    "prism:publicationName": "2013 International Conference on Control, Computing, C",  
    "prism:isbn": [...]  
},  
    "prism:pageRange": null,  
    "prism:coverDate": "2013-12-31",  
    "prism:coverDisplayDate": "2013",  
    "prism:doi": "10.1109/ICCCCM.2013.6648897",  
    "citedby-count": "0",  
    "prism:aggregationType": "Conference Proceeding",  
    "subtype": "ed",  
    "subtypeDescription": "Editorial",  
    "openaccess": "0",  
    "openaccessFlag": false,  
    "author": [...]  
},  
    "source-id": "21100276776"
```

A vertical list of years from 2013 to 2023, each preceded by a yellow folder icon. The years are: 2023, 2022, 2021, 2020, 2019, 2018, 2017, 2016, 2015, 2014, and 2013.

DATA COLLECTION



- Data Scraped
 - Author Details

A screenshot of a Scopus author profile page for Meimoun, Julie. The page shows basic metrics: 169 citations, 9 documents, and an h-index of 6. It also includes links to edit the profile, connect to ORCID, and view more metrics. The URL in the browser bar is https://www.scopus.com/authid/detail.uri?authorId=57200608089.

Explore this author profile on Scopus Preview
View limited highlights of a Scopus-generated author profile with Scopus Preview. To view the complete profile, check access through your organization. [Learn more about Scopus profiles.](#)

Check access

Meimoun, Julie

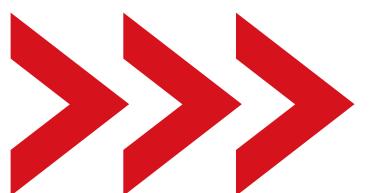
Université de Lille, Lille, France | 57200608089 | Connect to ORCID | Is this you? Connect to Mendeley account

169 Citations by 162 documents | 9 Documents | 6 h-index View h-graph | View more metrics >

Edit profile | More

9 Documents | Impact | Cited by 162 documents | 0 Preprints | 37 Co-Authors | 0 Topics | 0 Awarded Grants

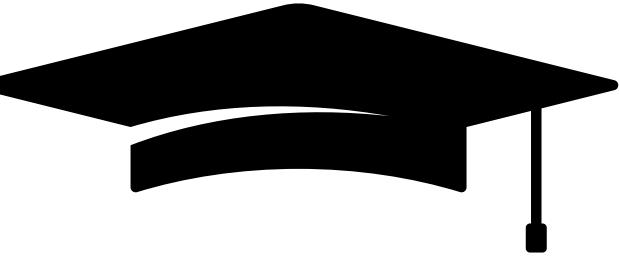
Browser Scraping



76182 records

```
{  
  "auid": "22980355000",  
  "author_stat": {  
    "citations": 25617,  
    "documents": 481,  
    "h_index": 80,  
    "co_authors_count": 11143,  
    "cited_by_count": 10306,  
    "preprints_count": 228  
  },  
  ...  
}
```

DATA COLLECTION



- Data Scraped
 - Author Details

```
author_id = author.get('@auid')

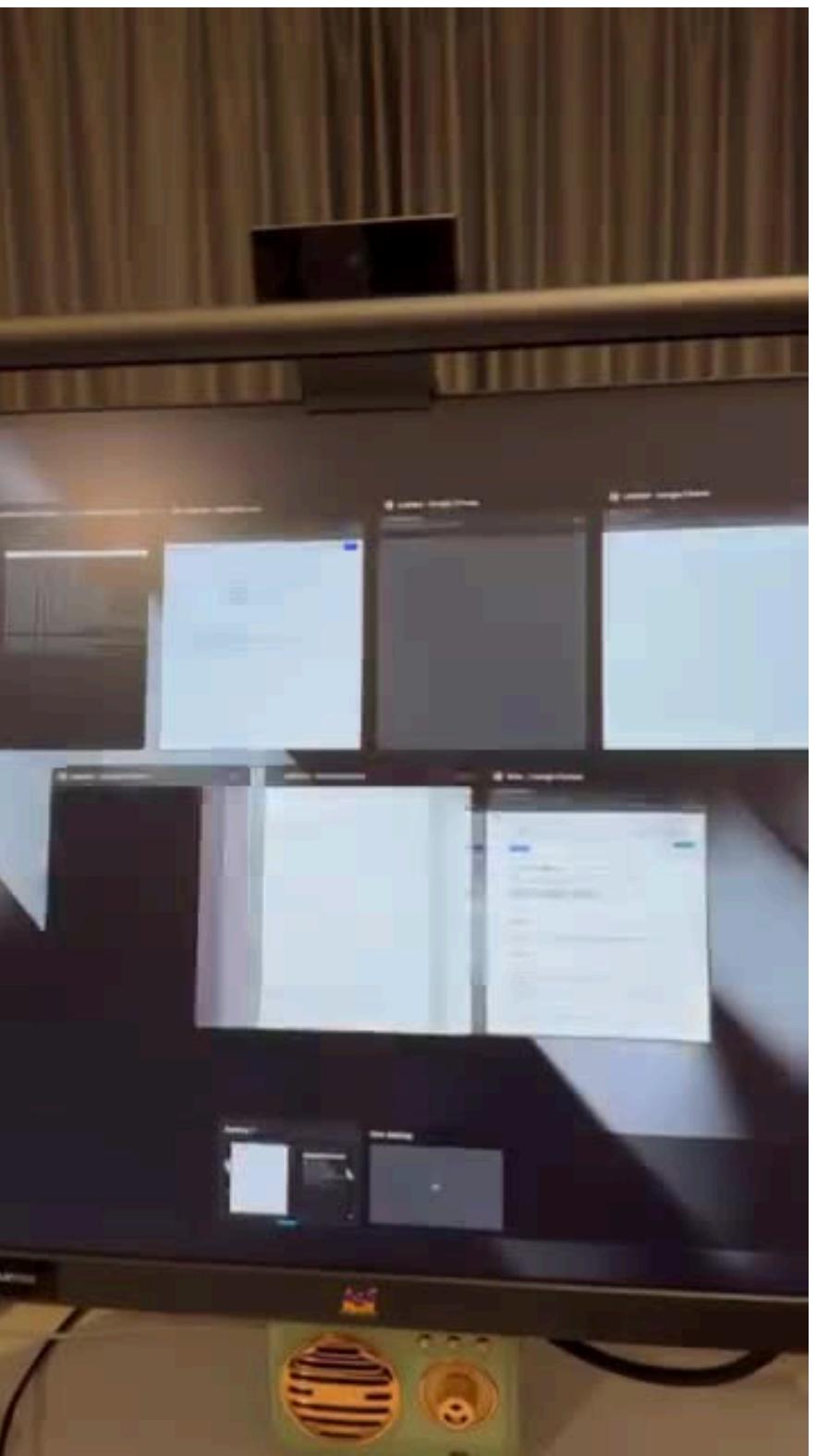
driver = self._create_webdriver()

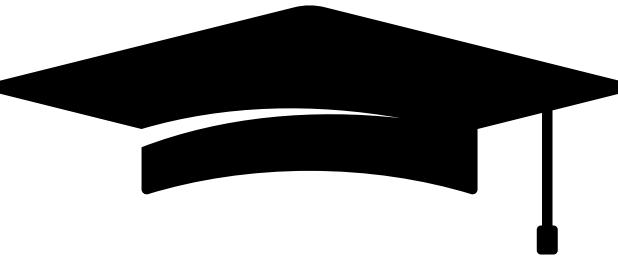
url = f'https://www.scopus.com/authid/detail.uri?authorId={author_id}'
driver.get(url)

return self._extract_metrics(driver, author_id)
```

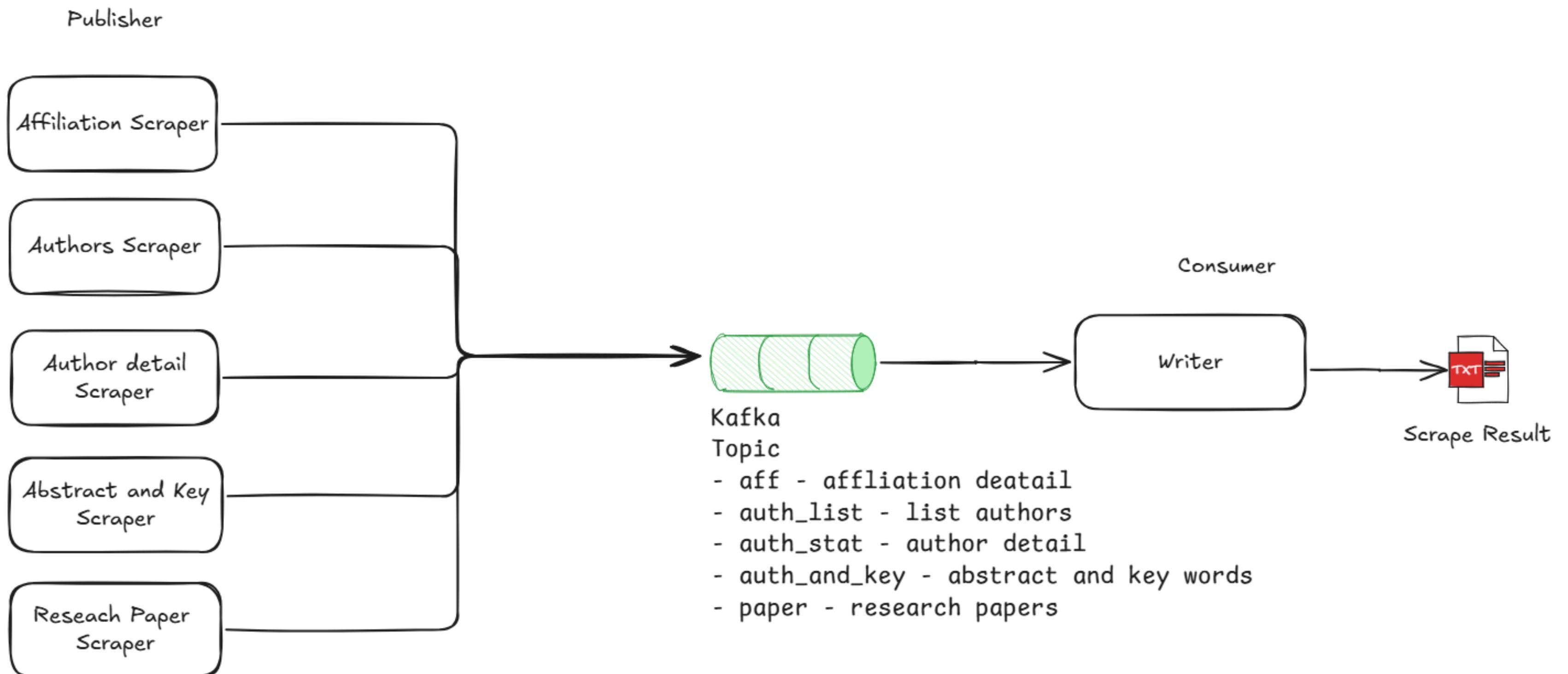


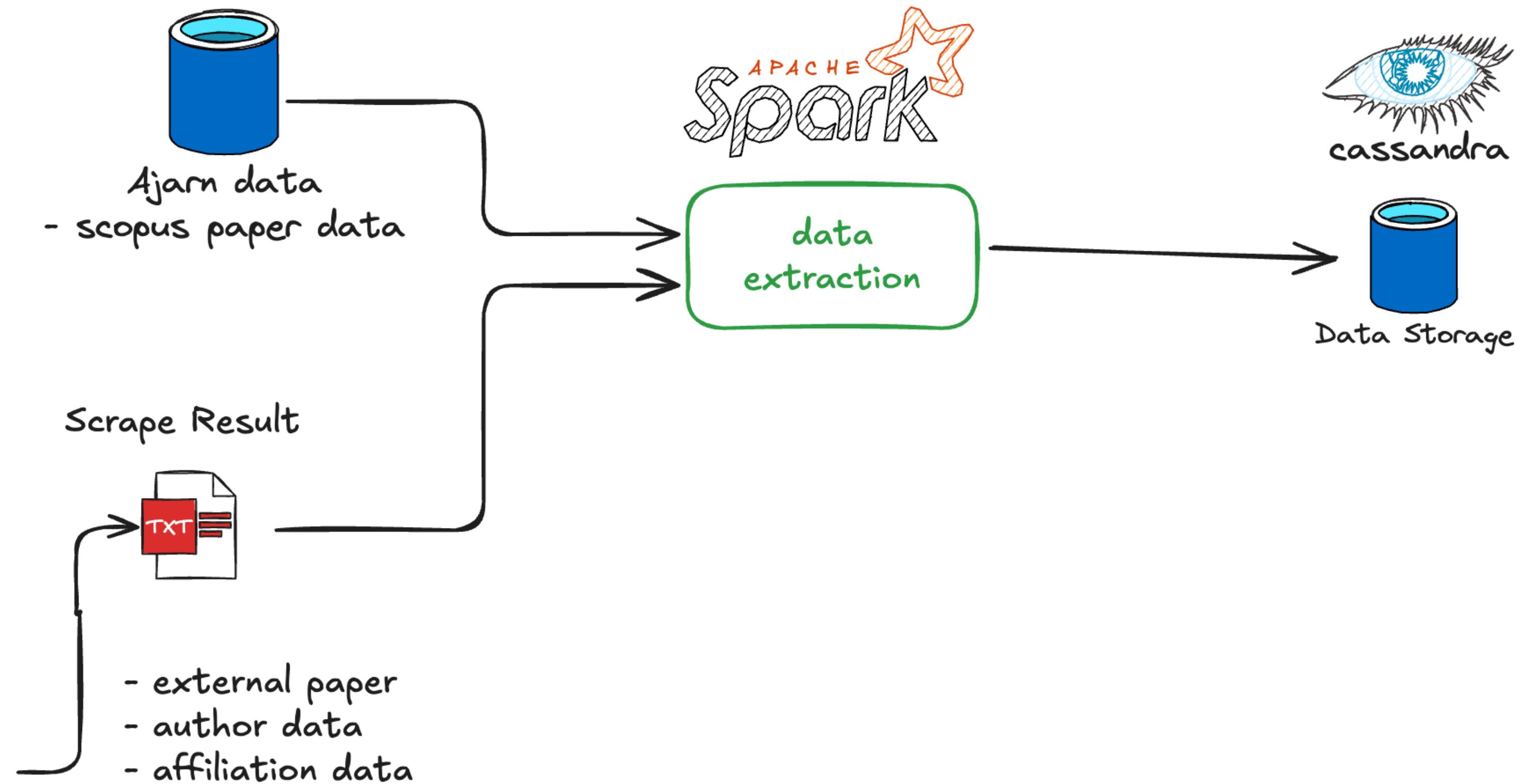
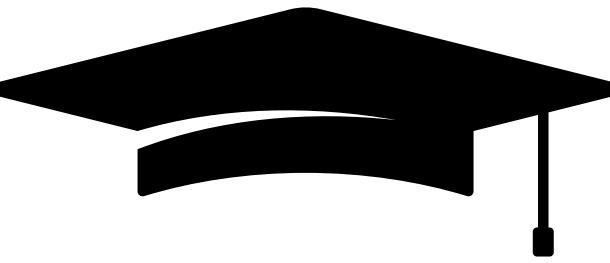
```
element = driver.find_element(By.CSS_SELECTOR, f"{selector} > span")
```





DATA INGESTION





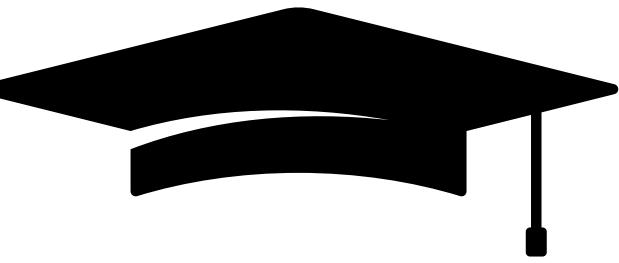
```
● session.execute("""  
    CREATE TABLE IF NOT EXISTS scopus_data.records (  
        doi text PRIMARY KEY,  
        title text,  
        abstract text,  
        document_type text,  
        source_type text,  
        publication_date text,  
        source_title text,  
        publisher text,  
        author_keywords text,  
        subject_code text,  
        subject_name text,  
        subject_abbrev text,  
        author_given_name text,  
        author_surname text,  
        author_url text,  
        author_affiliation text,  
        author_details text,  
        funding_details text,  
        ref_count int,  
        open_access boolean,  
        affiliation text,  
        language text,  
        cited_by int,  
        status_state text,  
        delivered_date text,  
        subject_area text  
    )  
""")
```

```
● session.execute("""  
    CREATE TABLE IF NOT EXISTS scopus_data.affiliations (  
        affiliation_id text PRIMARY KEY,  
        preferred_name text,  
        name_variants list<text>,  
        documents_count int,  
        authors_count int,  
        country text,  
        city text,  
        state text,  
        street_address text,  
        postal_code text,  
        contact_url text,  
        hierarchy_ids list<text>  
    )  
""")
```

```
● session.execute("""  
    CREATE TABLE IF NOT EXISTS scopus_data.author_metrics (  
        author_id text PRIMARY KEY,  
        citations bigint,  
        documents int,  
        h_index int,  
        co_authors_count int,  
        cited_by_count bigint,  
        preprints_count int  
    )  
""")
```

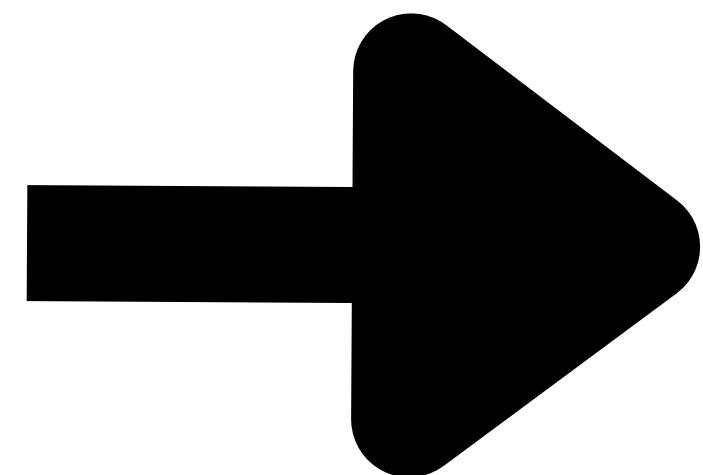
- records 62781
- author_metrics: 76182
- affiliations: 7625

Data Processing



- records 62781

```
 1 session.execute("")  
 2     CREATE TABLE IF NOT EXISTS scopus_data.records (  
 3         doi text PRIMARY KEY,  
 4         title text,  
 5         abstract text,  
 6         document_type text,  
 7         source_type text,  
 8         publication_date text,  
 9         source_title text,  
10         publisher text,  
11         author_keywords text,  
12         subject_code text,  
13         subject_name text,  
14         subject_abbrev text,  
15         author_given_name text,  
16         author_surname text,  
17         author_url text,  
18         author_affiliation text,  
19         author_details text,  
20         funding_details text,  
21         ref_count int,  
22         open_access boolean,  
23         affiliation text,  
24         language text,  
25         cited_by int,  
26         status_state text,  
27         delivered_date text,  
28         subject_area text  
29     )  
30     "")
```

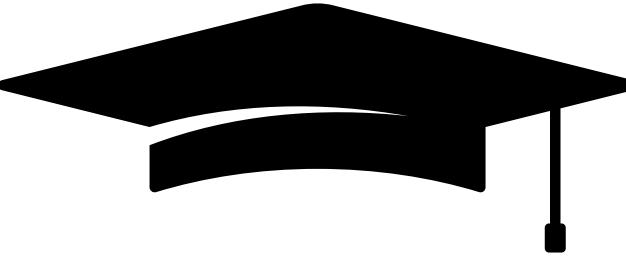


Abstract

Title

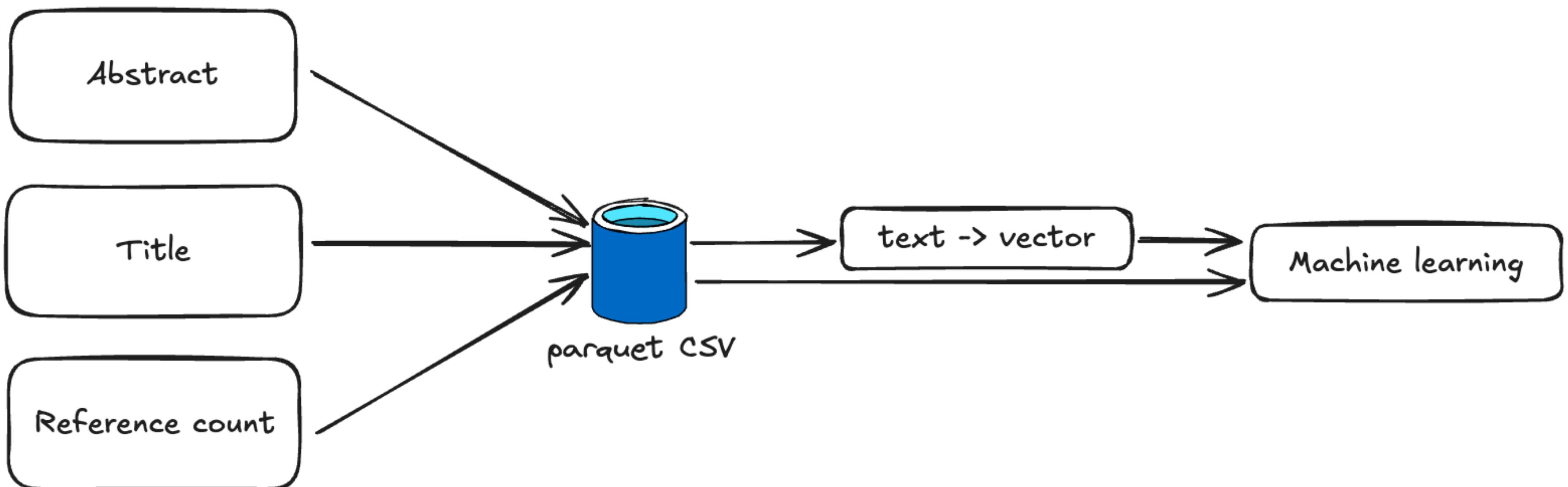
Reference count

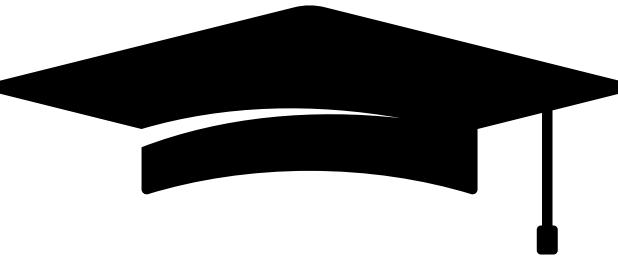
Data Processing



features selections

- reduce features following feature important model
- for more explainable and faster runtime





Data Processing

text feature extraction (bag of words)

- count how frequency that words appear in datas
- normalize by total of that words in documents

$\text{tf}(t, d)$

	blue	bright	can	see	shining	sky	sun	today
1	1/2	0	0	0	0	1/2	0	0
2	0	1/3	0	0	0	0	1/3	1/3
3	0	1/3	0	0	0	1/3	1/3	0
4	0	1/6	1/6	1/6	1/6	0	1/3	0

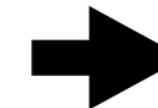
X

$\text{idf}(t, D)$

	blue	bright	can	see	shining	sky	sun	today
	0.602	0.125	0.602	0.602	0.602	0.301	0.125	0.602

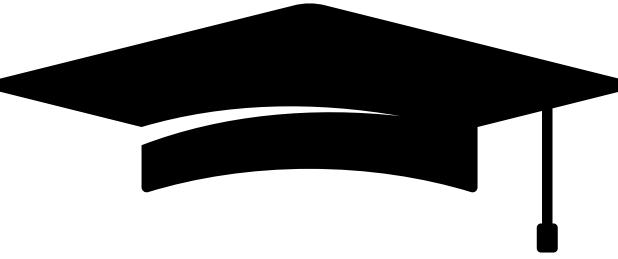
$$\text{tfidf}(t, d, D) = \text{tf}(t, d) \cdot \text{idf}(t, D)$$

- TF-IDF: Multiply TF and IDF scores, use to rank importance of words within documents



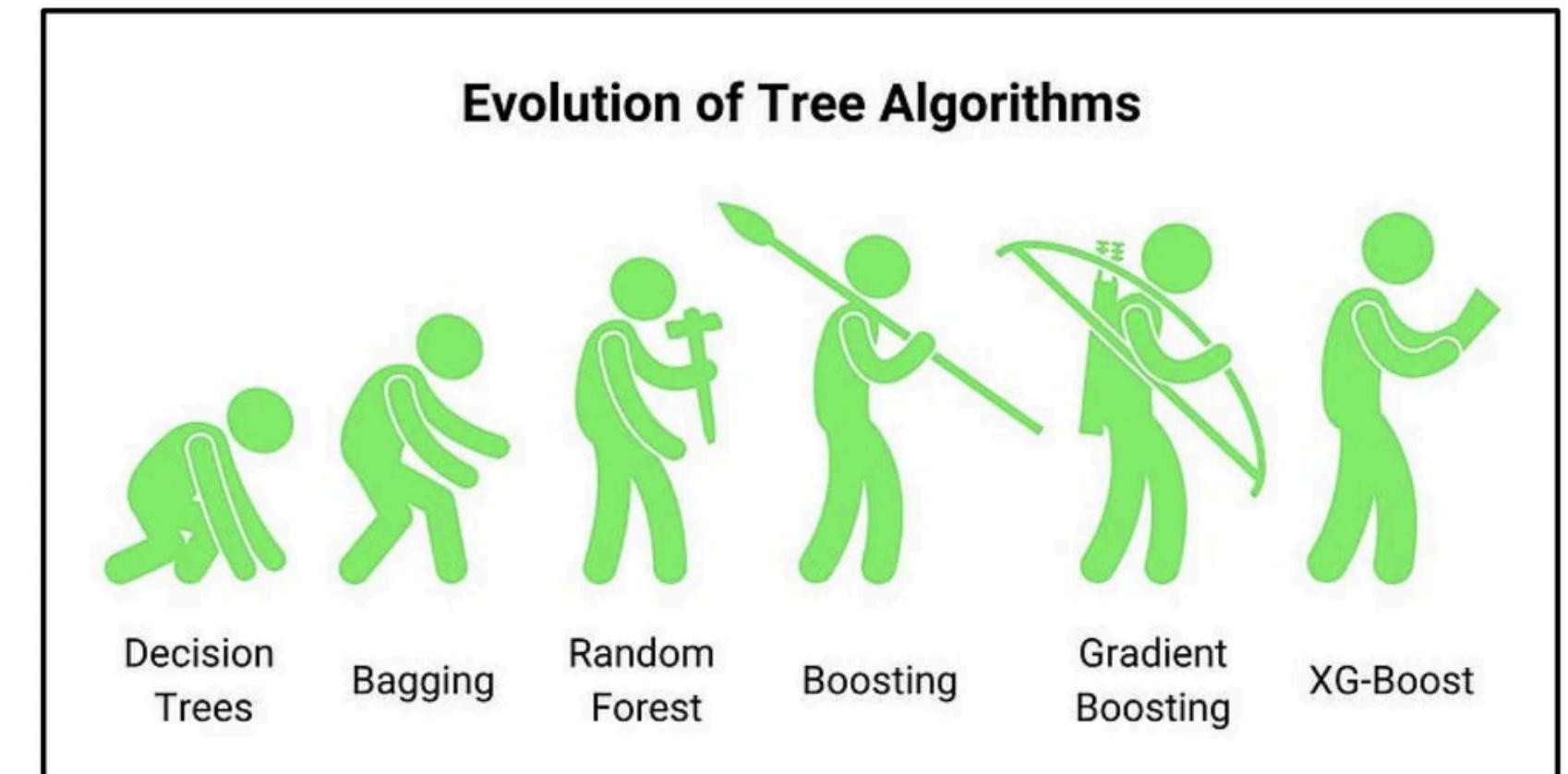
- Most important word for each document is highlighted

	blue	bright	can	see	shining	sky	sun	today
1	0.301	0	0	0	0	0.151	0	0
2	0	0.0417	0	0	0	0	0.0417	0.201
3	0	0.0417	0	0	0	0	0.100	0.0417
4	0	0.0209	0.100	0.100	0.100	0	0.0417	0



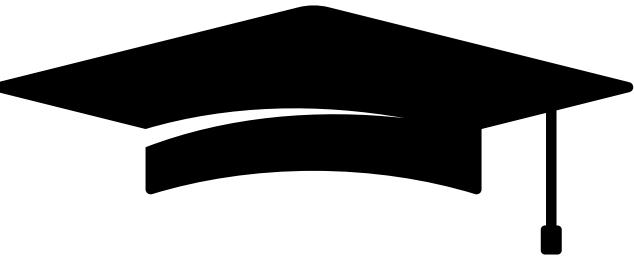
Model Selections

- classification -> probability
- tree based models
 - **XGBoost**
 - CatBoost
 - LightGBM
- hyperparameters

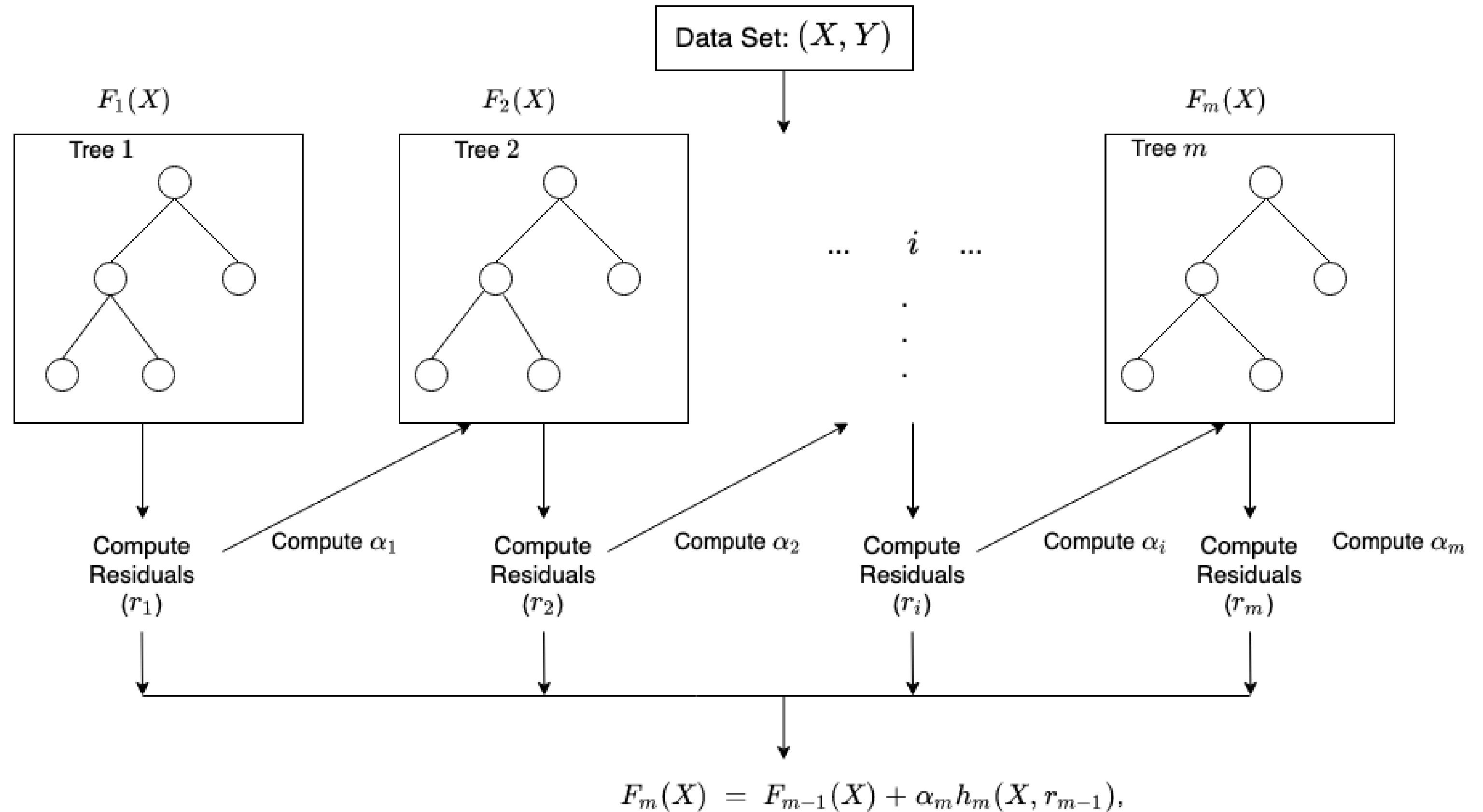


AutoGluon

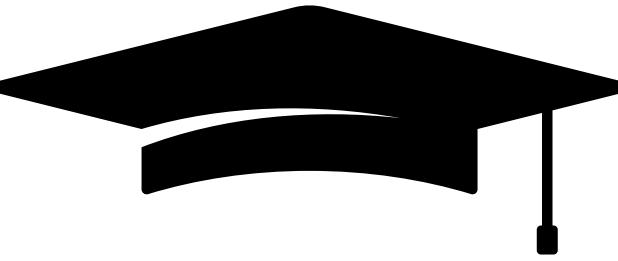
explainable model



XGBoost

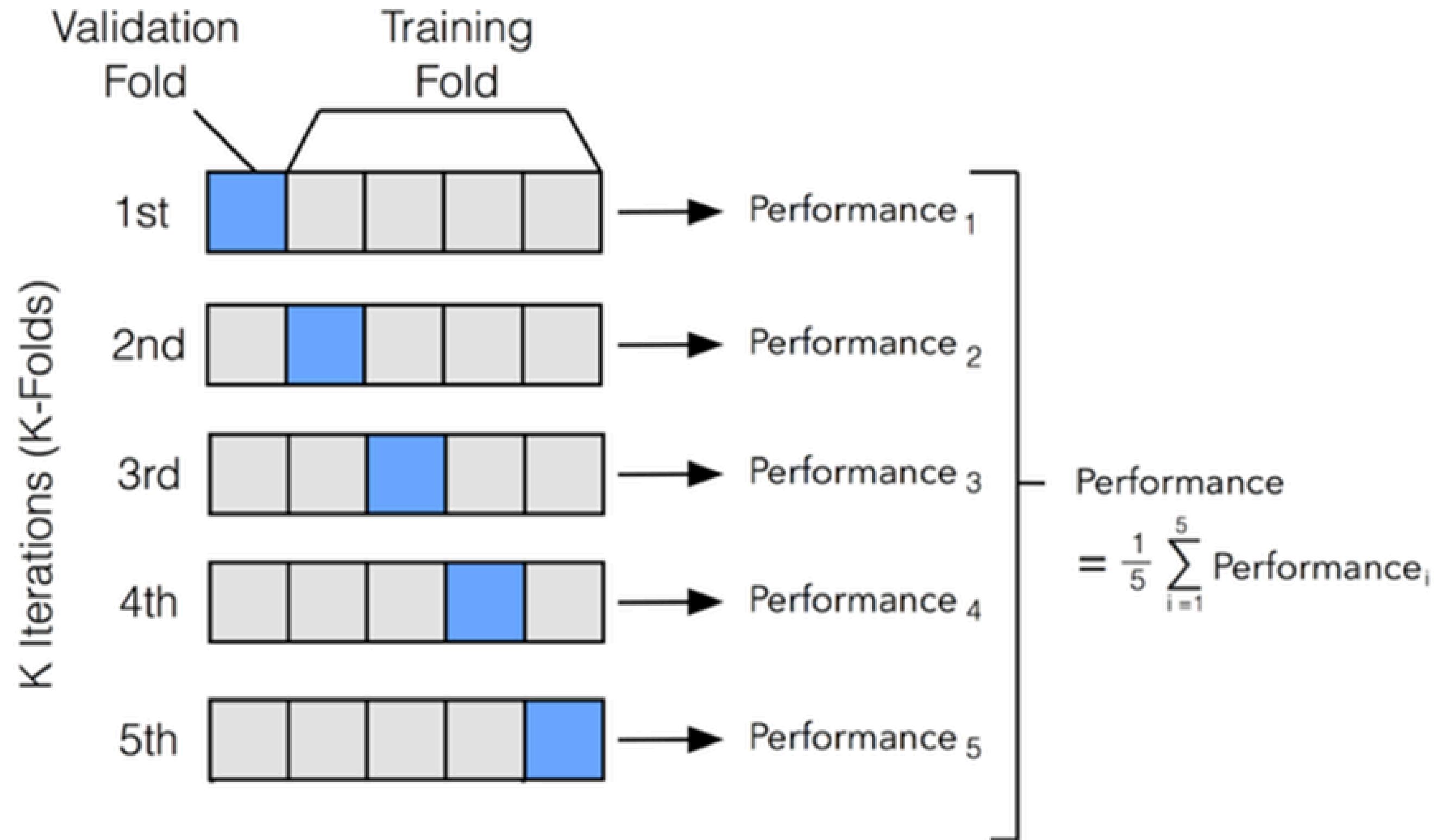


where α_i , and r_i are the regularization parameters and residuals computed with the i^{th} tree respectively, and h_i



training method

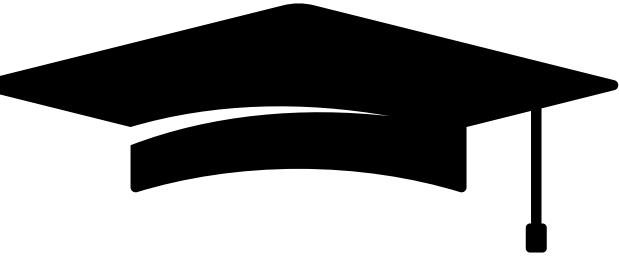
- **split datasets to k-fold then train k models**



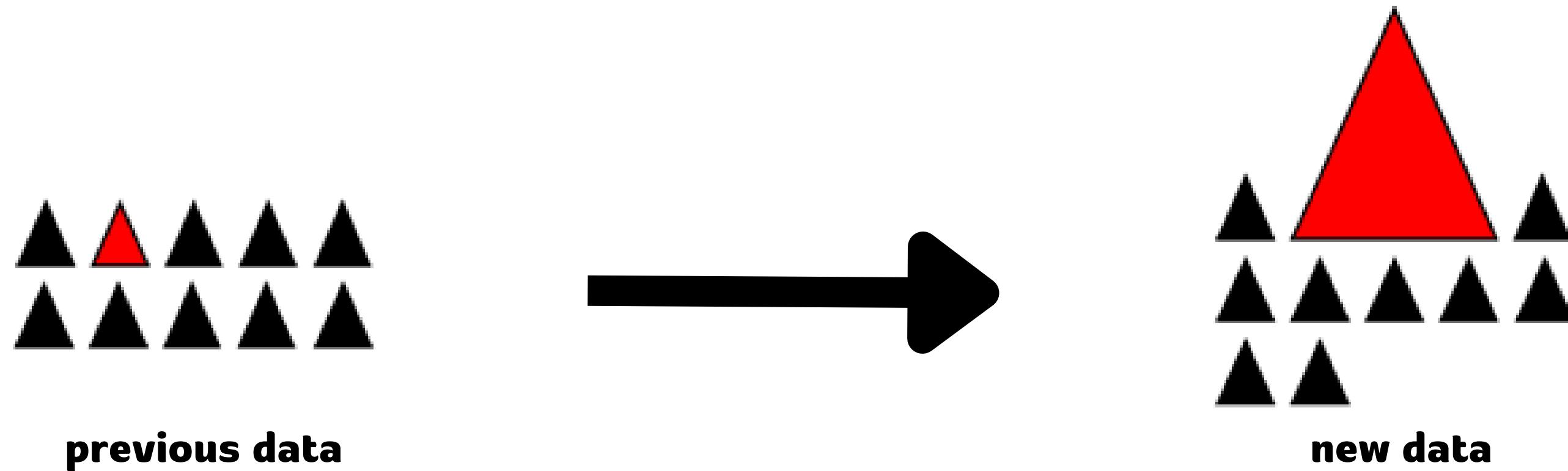
bag of words

not good enough

dealing with imbalance



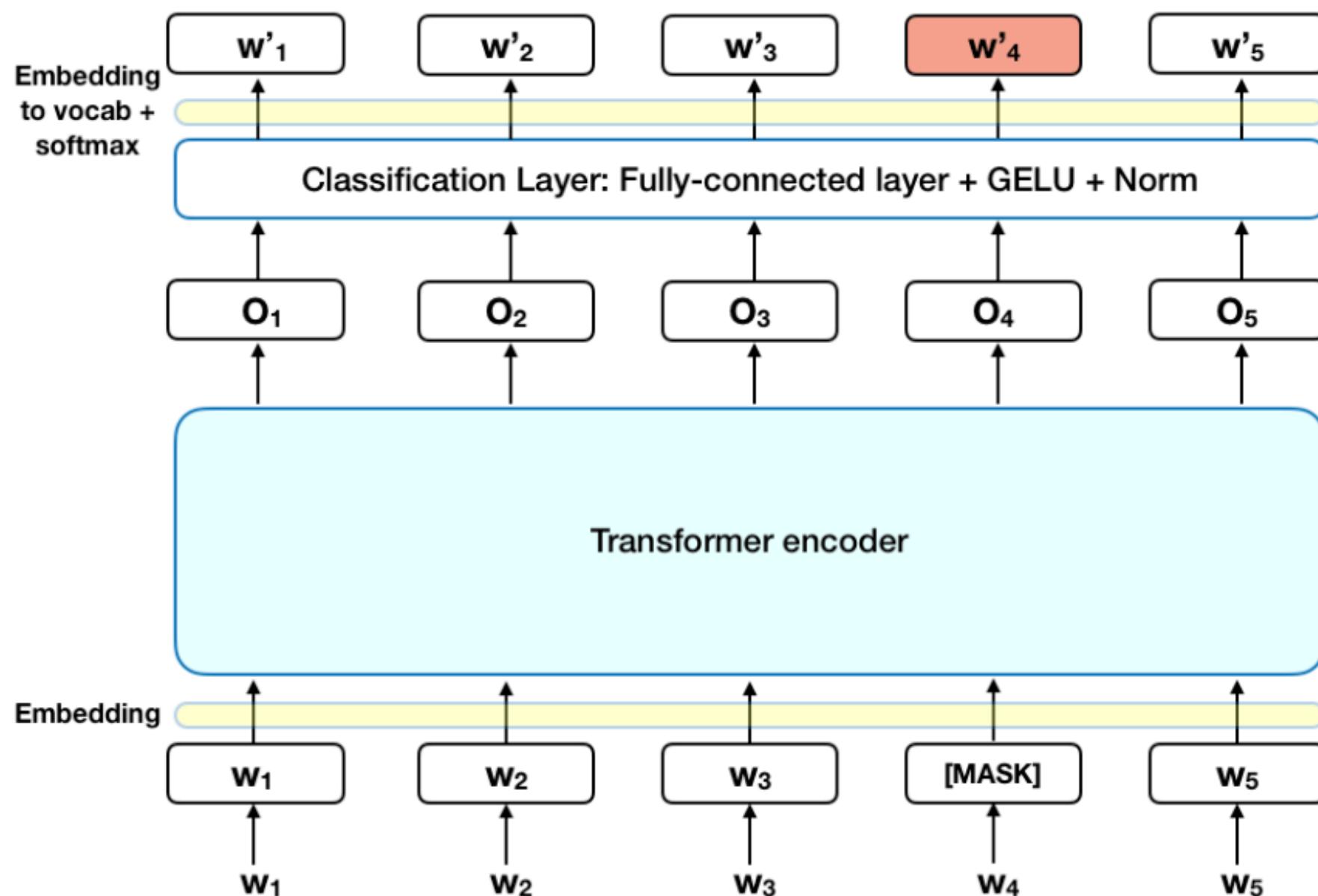
- **weight loss following distribution of class**
- have true class ~7 %



Data Processing

text feature extraction (Deep learning)

- use Bert that train with paper data (SciBert)
- have special token for papers
- 512 + 512 tokens



allenai/scibert_scivocab_uncased

SCIBERT: A Pretrained Language Model for Scientific Text

Iz Beltagy Kyle Lo Arman Cohan

Allen Institute for Artificial Intelligence, Seattle, WA, USA

{beltagy, kylel, armanc}@allenai.org

model results

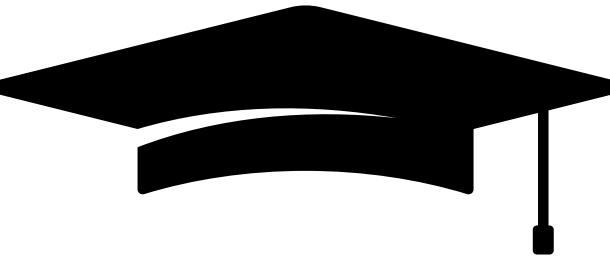


Overall Model Accuracy

94.3%

Precision, Recall, and F1-Score by Class





model results

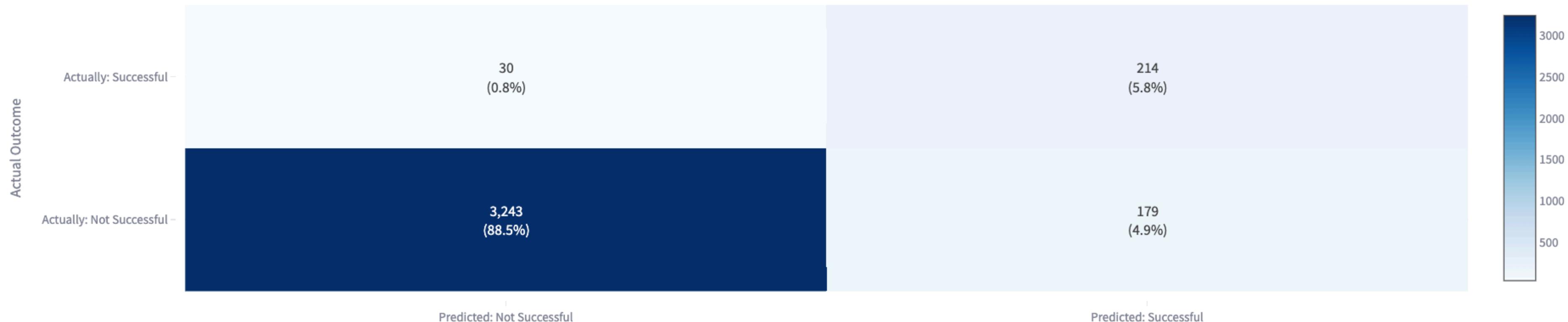
Recall:

- Correctly identifies **87.7%** of actually successful scholars

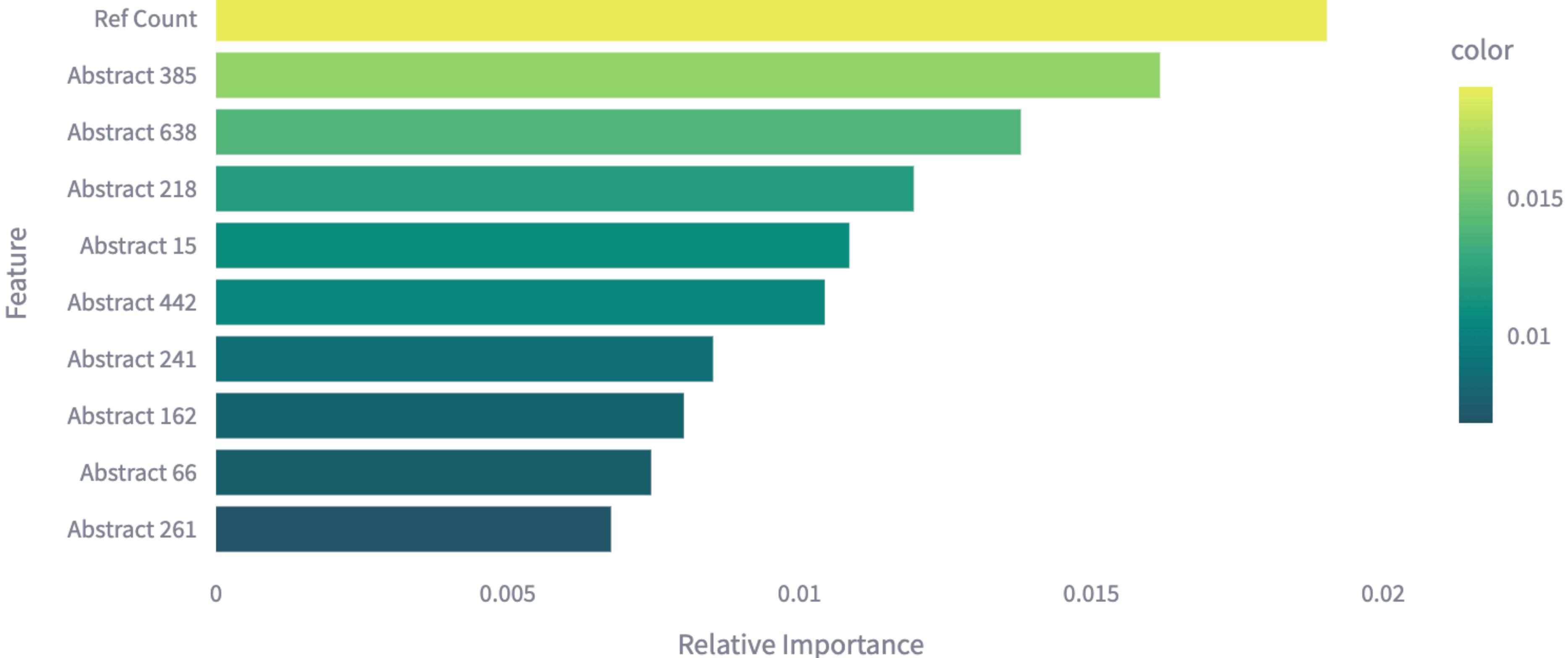
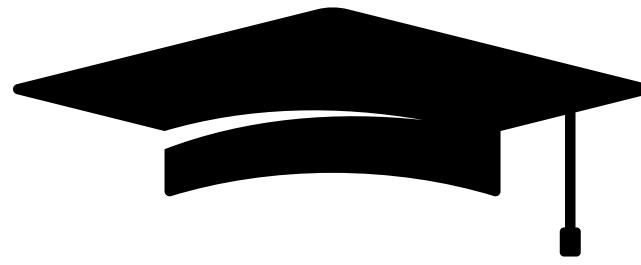
Dataset Distribution:

- Total scholars: 3666
- Successful: 244 (**6.7%**)
- Non-successful: 3422

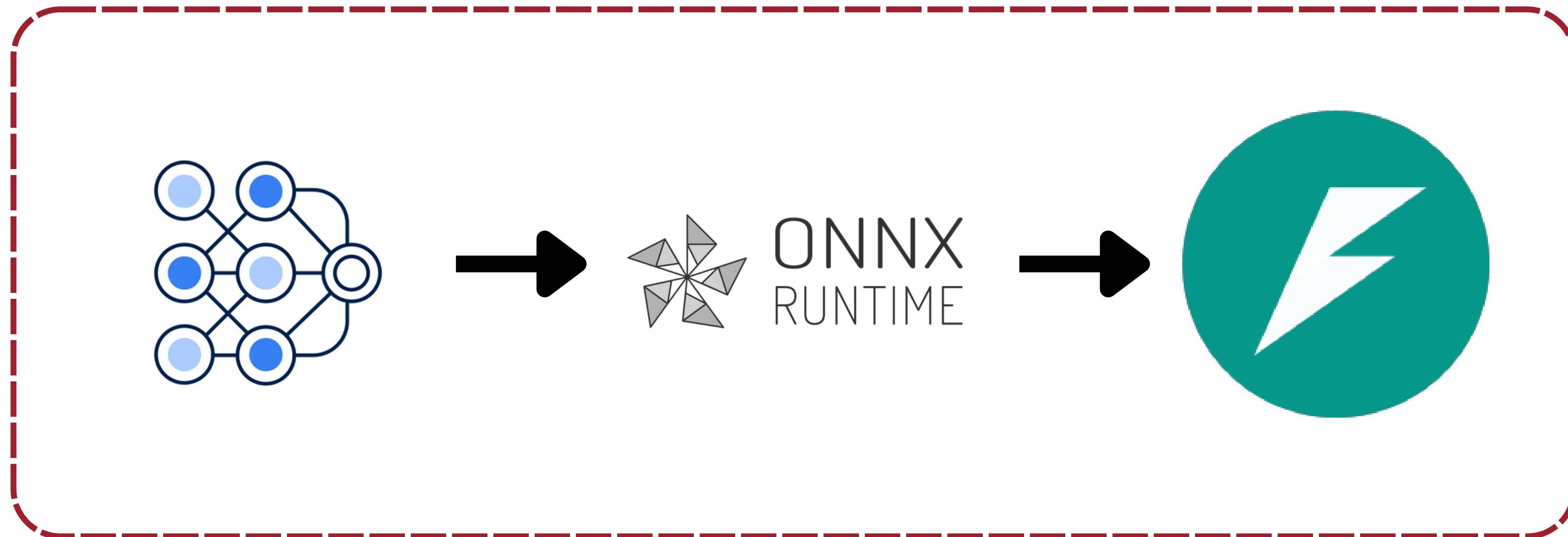
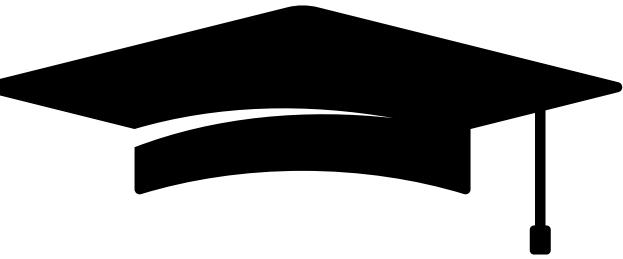
Confusion Matrix (Count and Percentage)



Top 10 Most Important Features for Prediction



deployment



DATA VISUALIZATION

Team Credential



SIRAWIT CHANABURANASAK

Software Developer



SADIT WONGPRAYON

Machine Learning engineer



PAPAWIN TANGCHITPORN

Software developer



SIRASIT TANAJIRAWAT

Project manager



SAKDIPAT SUKHANEKUL

Software developer