**Big Scientific Data and Text Analytics group :
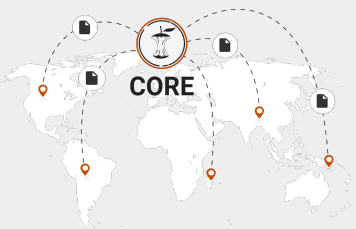AI for open and responsible research**

CORE

2023

# Applications of AI in Academic Libraries and Archives:

**Machine learning from and for open research.**

Prof. Petr Knoth

The Open University

# Big Scientific Data and Text Analytics group : AI for open and responsible research

**CORE** is the world's most used aggregator of **Open Access** papers, collating and enriching content from over **11,000 repositories**.

- **>20 Million** monthly active users (MAU)
- **34 Million** full-text research papers hosted by CORE.
- **260 Million** metadata records

Signatory of Principles of Open Scholarly Infrastructure **(POSI)**

CORE delivers **services** for HEIs, researchers, funders and commercial partners, offering seamless access to research.

| Content discovery | Raw data services | Managing content |
|---|---|---|
| Search | API | Repository Dashboard |
| Recommender | Dataset | Identifiers |
| Discovery | FastSync | OAI Resolver |

Providing seamless access to open research for humans and machines.

## Commercial Partners

ontochem IT SOLUTIONS

CAKTUS

cypris

EZASSI — Technology to innovate faster.

turnitin

- **Innovation and trends analysis**
- **Plagiarism detection**
- **Fact checking**
- **Finance**
- **Health**

## Institutional Members

UNIVERSITY OF OXFORD

Lancaster University

University of Huddersfield — Inspiring global professionals

University of BRISTOL

QUEEN'S UNIVERSITY BELFAST

THE UNIVERSITY OF CHICAGO

University of St Andrews FOUNDED 1413

a.r.u.

UNIVERSITY OF BIRMINGHAM

White Rose university consortium

UNIVERSITY of York

UNIVERSITY OF LEEDS

University of Sheffield

32 supporting or sustaining members

## Research areas

→ AI Applications in Research Evaluation (e.g. citation type classification, bibliometrics, impact assessment)

→ Automatic Expert Finder systems (e.g. for peer-review and grant applications)

→ Deduplication, document classification, rapid systematic reviews

→ Research graphs: entity extraction (affiliation, author, etc.)

→ Research recommender systems and academic search

Dr. Petr Knoth : Senior Research Fellow in Text and Data Mining petr.knoth@open.ac.uk

The Open University

# CORE and the OA landscape

**CORE's mission is**

to index all open access research worldwide and deliver unrestricted access for all.

We are here to support and advance the Open Access / Open Research movement

**WE ARE**

the world's **most used** collection of open access research papers from repositories

**WE ARE**

a **not-for-profit** scholarly infrastructure dedicated to the open access mission, **adopters of POSI** principles.

**WE**

**provide solutions** for content management, discovery and scalable machine access to research.

**WE**

**serve the global network** of repositories and journals by increasing discoverability and reuse of open access content.

# CORE - An adopter of POSI

The Principles of Open Scholarly Infrastructure

- ⤍ **Governance**
- ⤍ **Sustainability**
- ⤍ **Insurance**

**CORE** is a mission-driven not-for-profit endeavour and a signatory of the **Principles of Open Scholarly Infrastructure.**

CORE

# CORE Community Governance



**Advisory Board**
- Advises on strategic directions
- Ensures mission alignment with the needs of the open research community

**Board of Supporters**
- Helping to identify requirements and prioritise the development roadmap
- Represents the interests of the global open repositories and journals network.

**The Open University Stakeholder Group**
- Assumes overall financial and legal responsibility for CORE's obligations.
- Provides institutional support and resources for CORE (e.g. HR, financial, legal, infrastructure).

**Research network representatives**
- Ensures relevance and provides guidance on effectively supporting the open research community.

**CORE Leadership & Management Team**
- Is responsible for the day to day operation of CORE.
- Takes operational decisions with guidance from the governance groups

CORE

# Outline

1. How can Artificial Intelligence and Machine Learning (**AI/ML**) applied to research papers benefit and **transform research**
2. The crucial role of repositories in providing **machine access** to research content.
3. Using AI/ML for **research intelligence** and improving repository workflows

CORE

# Outline

The aggregation of repository content can offer the foundation for a whole host of text mining activities to be developed on top of the content. Text and data mining are becoming valuable analytical methods that allow researcher to discover interesting patterns and extract new knowledge from a corpus of content. Repository collections contain all kinds contain rich information, which could be further used, combined and analyzed through text mining techniques. A growing number of services are being developed to support these types of service.[30] As text and data mining techniques in research are more widely adopted, repositories and the broader community will need to

# How can AI/ML transform research

# The importance of open research literature

Research literature documents the knowledge we have assembled as human species.

# The wide variety of use cases over research literature

- A limited number of use cases can be satisfied with a sample of scholarly content.

- Many use cases require machine access to all existing research from everywhere and always up-to-date.

- High cost when a repository does not participate in the open network, by not providing machine access. Some use cases significantly affected.

**Undefined sample use cases**

**Direct use cases** (answer specific question, gain insights)

**Indirect use cases** (information access/organisation)

- Detecting plagiarism in newly submitted publications
- Semantometrics
- Matching publications to suitable reviewers
- Analysing the growth of science and research trends

- Recommending papers, collaborators, venues, etc.
- Identifying different versions of the same publication (e.g. pre-print vs post-print)

- Detecting if a publication is within scope of a venue
- Detecting importance, sentiment, or type of citations

- Summarising research findings
- Building citation indices
- Extracting interesting words and phrases
- Categorising publications into fields of study
- Detecting publication type
- Extracting citations for use in bibliometrics

**A priori defined sample use cases**

# AI for systematic reviews

➔ **Systematic reviews**

 ◆ Time consuming

➔ **Rapid reviews**

 ◆ Limitation on the number of outcomes, interventions and comparators

➔ **Living reviews**

 ◆ Live updates to historic systematic reviews with the help of recommender system

# AI for systematic reviews

→ Involves many steps

→ Some of teh most-time consuming can be automated

| Step/Task | Description | Stage |
|---|---|---|
| 1. Formulate review question | Decide on the research question of the review | Preparation |
| 2. Find previous systematic reviews | Search for SR that answers the same question, (part of scoping the literature in EFSA guidance) | Preparation |
| 3. Write the protocol | Provide an objective, reproducible, sound methodology for peer review | Write up |
| 4. Devise search strategy | Decide on databases and keywords to find all relevant trials | Preparation |
| 5. Search | Aim to find all relevant citations even if many irrelevant ones are included | Retrieval |
| 6. De-duplicate | Remove identical citations | Retrieval |
| 7. Screen abstracts | Based on titles and abstracts, remove definitely irrelevant trials | Screening |
| 8. Obtain full text | Download or request copies from authors | Retrieval |
| 9. Screen full text | Exclude irrelevant trials | Screening |
| 10. Snowball | Follow citations from included trials to find additional ones | Retrieval |
| 11. Extract data | Extract relevant information (either quantitative or qualitative) to help with the synthesis and conclusions | Synthesis |
| 12. Critical appraisal | Assessing the risk of bias in the included studies | Critical Appraisal/ Synthesis |
| 13. Synthesize data | Convert extracted data to a common representation considering the results from the critical appraisal (if /when applicable) | Synthesis |
| 14. Re-check literature | Repeat search to find new literature published since the initial search | Retrieval |
| 15. Meta analyse | Statistically combine the result from all included trials | Synthesis |
| 16. Write up review | Produce and publish final report | Write up |

CORE

# AI for systematic reviews

Documents search → 2072 papers | 302 papers (14%) | 42 papers (2%) → Data extraction

*Title and abstract screening*

*Full text screening*

CORE

# AI for systematic reviews



CORE

# AI for systematic reviews



Which neural retrieval models exist for domain specific search? `Screen papers`

Description: Neural ranking and retrieval models show great performance gains in a setting where the models are trained with a large number of labeled training data (MS Marco). However it is not clear how these findings There can be a different task setting which requires a different modeling approach of the neural model of the problem and there can be the lack of training data. I want to know more about existing approaches of neural ret retrieval settings.

Search queries: domain specific Neural retrieval models   Deep learning for domain specific information retrieval   Dense retrieval BERT for domain specific retrieval   Transformer-based domain specific retrieval models

Inclusion criteria: Paper about neural retrieval model   Paper about BERT like model   Paper using Transformer
Exclusion criteria: Paper written in language other than English   Paper about web search   Paper about statistical retrieval model   Paper older than 2014

Show review statistics •

| # | Title | Authors | Year | Journal | Citations | URL | PDF | Screened |
|---|-------|---------|------|---------|-----------|-----|-----|----------|
| 1 | Neural IR for Domain-Specific Tasks | Oscar Espitia, G. Pasi | | IIR | – | 🌐 | — | Yes, edit |
| 2 | Using Siamese Graph Neural Networks for Similarity-Based Retrieval in Process-Oriented Case-Based Reasoning | Maximilian Hoffmann, Lukas Malburg, P. Klein, R. Bergmann | | ICCBR | 7 | 🌐 | — | Yes, edit |
| 3 | Image Retrieval by Fusion of Features from Pre-trained Deep Convolution Neural Networks | Vijayakumar Bhandi, K. S. Sumithra Devi | 2019 | 2019 1st International Conference on Advanced Technologies in Intelligent Control, Environment, Computing & Communication Engineering (ICATIECE) | 6 | 🌐 | | Yes, edit |
| 4 | InPars: Data Augmentation for Information Retrieval using Large Language Models | L. Bonifacio, H. Abonizio, Marzieh Fadaee, Rodrigo Nogueira | 2022 | ArXiv | 3 | 🌐 | 📕 | Yes, edit |
| 5 | Re-ranking Biomedical Literature for Precision Medicine with Pre-trained Neural Models | Jiazhao Li, Adharsh Murali, Q. Mei, V.G.Vinod Vydiswaran | 2020 | 2020 IEEE International Conference on Healthcare Informatics (ICHI) | – | 🌐 | — | Yes, edit |
| 6 | Rapid Probabilistic Interest Learning from Domain-Specific Pairwise Image Comparisons | Michael Burke, Siyabonga Mbonambi, Purity Molala, R. Sefala | 2017 | | 1 | 🌐 | — | Yes, edit |
| 7 | Word Embedding Models for Query Expansion in Answer Passage Retrieval | Nirmal Roy | | | 1 | 🌐 | — | Yes, edit |
| 8 | Ranking Model for Domain Specific Search | Priyanka Jadhav, Vaishali S. Pawar, C. Jadhav, Nidhi R. Sharma | | | – | 🌐 | — | Yes, edit |
| 9 | Improving Passage Retrieval with Zero-Shot Question Generation | Devendra Singh Sachan, M. Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, J. Pineau, Luke Zettlemoyer | 2022 | ArXiv | 4 | 🌐 | 📕 | Yes, edit |
| 10 | Convolutional Neural Network Based use Surveillance Videos for Recognizing Human Actions Based on Machine Learning | Dipak Daitkar, Divyesh Patil, Akshay Desai, Prasad Kawade, Prof. M. R. Bendre | | | – | 🌐 | — | Yes, edit |

CORE

# AI for systematic reviews

## Neural Passage Retrieval with Improved Negative Contrast 🌐

**Abstract:** In this paper we explore the effects of negative sampling in dual encoder models used to retrieve passages for automatic question answering. We explore four negative sampling strategies that complement the straightforward random sampling of negatives, typically used to train dual encoder models. Out... Show full abstract

*Jing Lu, Gustavo Hernández Ábrego, Ji Ma, Jianmo Ni, Yinfei Yang*

2020 — ArXiv

### 1. Relevance *

| Domain relevance | very relevant | somewhat relevant | not relevant |
|---|---|---|---|

| Topic relevance | very relevant | somewhat relevant | not relevant |
|---|---|---|---|

### 2. Inclusion criteria

| Paper about neural retrieval model | Yes | Not sure | No |
|---|---|---|---|
| Paper introducing NEW retrieval model | Yes | Not sure | No |
| Paper about BERT like model | Yes | Not sure | No |
| Paper using Transformer | Yes | Not sure | No |

### 3. Exclusion criteria *

| Paper written in language other than English | Yes | Not sure | No |
|---|---|---|---|
| Paper not about domain specific search | Yes | Not sure | No |
| Paper about statistical retrieval model | Yes | Not sure | No |
| Paper older than 2014 | Yes | Not sure | No |
| Only title is available | Yes | Not sure | No |

### 4. Descriptive reason

### 5. Decision based on title and abstract *

| Include | Not sure | Exclude |
|---|---|---|

### 6. Did you know this paper before? *

| I knew and read the full paper before | I knew the paper but not read the full paper before | I did not know it before |
|---|---|---|

### 7. Did you know any of the authors before? *

| Yes, I knew at least one of the authors | No, I did not know any of the authors |
|---|---|

CORE

# AI for citation typing and research assessment

11 of 34 **REF2014** Peer Review Panels used citation data to 'inform' their decisions

**REF GPA** results highly correlated with citation data in these domains

Addition of citation type information can allow for better modelling of how research is being used.

Potential for development of new metrics that leverage enhanced citation information

*'The pilot exercise concluded that citation information is not sufficiently robust to be used formulaically or as a primary indicator of quality in the REF'*

HEFCE. Report on the pilot exercise to develop bibliometric indicators for the Research Excellence Framework;

| | UoA | mn2017 | med2017 | mn2014 | med2014 |
|---|---|---|---|---|---|
| 1 | Chemistry | 0.663 | 0.802 | 0.637 | 0.738 |
| 2 | Biological Sciences | 0.782 | 0.797 | 0.688 | 0.785 |
| 3 | Aero. Mech. Chem. Engineering | 0.771 | 0.758 | 0.745 | 0.760 |
| 4 | Social Work and Policy | 0.697 | 0.752 | 0.629 | 0.635 |
| 5 | Computer Science and Informatics | 0.715 | 0.743 | 0.720 | 0.678 |

CORE

# AI for citation typing and research assessment

- Knowing not only that something was cited, but WHY it was cited.

- Built ACT Dataset of >11,000 citations annotated by authors according to classification schema

- Ran 2 Shared Tasks to establish benchmarks for SoA classification models using ACT and extended ACT2 datasets

- Currently investigating extended / dynamic citation contexts to improve model performance

| Citation Function | Examples |
| --- | --- |
| BACKGROUND | Most of the participatory models to design educational games are founded on educational theories and game design (see for example: Amory, 2007; #CITATION_TAG). |
| COMPARES_CONTRASTS | Similar observations have been made in the past [30] [31] [32] [33] [34], although others have reported either no relationship or a negative association with SES [#CITATION_TAG]. |
| EXTENSION | This database is the result of a mandatory questionnaire about the home to work displacements and the mobility management measures at large workplaces in Belgium (#CITATION_TAG). |
| FUTURE | We are thus exploring the option of using datasets such as CrossRef 12, Dimensions 13, OpenCitations [11], and Core [#CITATION_TAG]. |
| MOTIVATION | To illustrate, consider the motivation given by #CITATION_TAG in developing their Bayesian account of word learning. |
| USES | The diffraction patterns from single crystal measurements were indexed with a home-made program based on the Fit2D software [#CITATION_TAG]. |

CORE

# A prototypical citation intent classification system

CORE

# Evaluation / shared tasks for citation classification

- Citation Context Classification (3C)
shared task

- ACT 2 dataset

- Same conditions for every team

Kunnath, Suchetha N.; Pride, David; Herrmannova, Drahomira and
Knoth, Petr **Overview of the 2021 SDP 3C Citation Context
Classification Shared Task**. In: *Proceedings of the Second
Workshop on Scholarly Document Processing, Association for
Computational Linguistics*, Stroudsburg, PA, pp. 150–158. URL:
https://aclanthology.org/2021.sdp-1.21/

(a) Subtask A

(b) Subtask B

Figure 2: Public Vs Private Macro F-Score performance on the leaderboard

CORE

Table 4. Institution-level Pearson correlations between machine learning predictions with 50% used for training and actual scores for articles 2014-18, following Strategy 1 (averaged across 10 iterations) and aggregated by institution for UoAs 1-11 and 16.

| UoA | Actual vs machine learning predicted average score | Actual vs machine learning predicted total score |
|---|---|---|
| 1:Clinical Medicine | 0.895 | 0.998 |
| 2:Public Health, Health Services and Primary Care | 0.906 | 0.995 |
| 3:Allied Health Professions, Dentistry, Nursing & Pharmacy | 0.747 | 0.982 |
| 4:Psychology, Psychiatry and Neuroscience | 0.844 | 0.995 |
| 5:Biological Sciences | 0.885 | 0.995 |
| 6:Agriculture, Food and Veterinary Sciences | 0.759 | 0.975 |
| 7:Earth Systems and Environmental Sciences | 0.840 | 0.986 |
| 8:Chemistry | 0.897 | 0.978 |
| 9:Physics | 0.855 | 0.989 |
| 10:Mathematical Sciences | 0.664 | 0.984 |
| 11:Computer Science and Informatics | 0.724 | 0.945 |
| 16:Economics and Econometrics | 0.862 | 0.974 |

# AI for citation typing and research assessment

Significant performance improvement of parameter updating methods across a variety of prompting strategies over promptless fine-tuning.

Dynamic context-based prompts significantly improve model scores for both datasets and surpass the performance on the 3C shared task bench-mark.

Kunnath, Suchetha N.; Pride, David and Knoth, Petr (2022) **Prompting Strategies for Citation Classification** In: *CIKM 2023*

# GPT Fact vs fiction

50 questions from across 10 different domains.

Each answer supported by 5 references to papers

72.5% of citations provided by GPT3.5 were fictional.
71.2% of citations provided by GPT4 were fictional.



| | | GPT3.5 | GPT4 |
|---|---|---|---|
| 🔴 | Fictional | 72.5% | 71.2% |
| 🟡 | Conflated | 5.5% | 9.2% |
| 🟢 | Factual | 22.0% | 19.6% |

Tell me a joke about ChatGPT responses not being based on scientific papers

Why did ChatGPT refuse to base its responses on scientific papers?

Because it didn't want to get lost in the references!

CORE

# ChatGPT



Tell me a joke about ChatGPT responses not being based on scientific papers

Why did ChatGPT refuse to base its responses on scientific papers?

Because it didn't want to get lost in the references!

# AI for credible trustworthy question answering (CORE-GPT)

**CORE** is the world's largest collection of Open Access papers, collating and enriching content from over **11,000** data providers.

**GPT** large language models*

- Can comprehend context and generate human-like text
- Can infer meaning from large-scale data



**CORE**

- **>20 Million** monthly active users
- **34 Million** full-text research papers hosted by CORE.
- **260 Million** metadata records



GPT-3

robotics → 1 → 2 → 3 → A

1- Convert word into vector    2- Magic    3- Convert vector into word

*Other large language models are available

@JayAlammar

# Introducing CORE-GPT

# CORE-GPT Results

# CORE-GPT Results



*What are the latest research trend in biomimicry?*

Biomimicry in architecture and building design is a promising research trend for reducing energy use and decarbonizing buildings, as it takes inspiration from nature to create more sustainable solutions [1]. Mechatronic and soft robotics, particularly in marine environments, are also taking inspiration from the animal kingdom to create high-performance robots, and there is significant research in energy provision, biomaterials, and locomotor designs [2]. Similarly, a recent biomimetic exploratory study aimed at designing a passive cooling unit for buildings in hot, arid regions by emulating biological cooling strategies in nature, such as cooling through animals' respiratory passages, thermal radiators, and airflow cooling in termite mounds [3]. Lastly, to address climate change and biodiversity loss, built environment-focused responses that take advantage of biomimicry could potentially make a contribution to solving these problems by interacting harmoniously with the living world [4].

1. https://core.ac.uk/works/131198542 – Biomimetic Buildings: Copying Nature for Energy Efficiency
2. https://core.ac.uk/works/8388662 – Research Trends and Future Perspectives in Marine Biomimicking Robotics
3. https://core.ac.uk/works/68990789 – Biomimetics for passive air conditioned design for buildings in the hot arid regions
4. https://core.ac.uk/works/24102079 – Ecosystem Services Analysis for the Design of Regenerative Urban Built Environments
5. https://core.ac.uk/works/123801488 – Holistic biomimicry: a biologically inspired approach to environmentally benign engineering

See more in CORE

CORE

# How does CORE-GPT work?



User query

...Generate an Elastic Search query for research papers, based on the question...

/v3/works/search
**API**

[work1, work2, work3, work4, work5]

Generate a comprehensive answer (but no more than **160** words) solely based on the content of the provided search results. Format the links to the papers as follows: {url:$url, abstract:$abstract}

CORE GPT validator

Results

Pride, David; Cancellieri, Matteo and Knoth, Petr (2022) **CORE-GPT: Combining Open Access research and large language models for credible, trustworthy question answering**. In: *TPDL 2023*

CORE

# How well does CORE-GPT work?



Pride, David; Cancellieri, Matteo and Knoth, Petr (2022) **CORE-GPT: Combining Open Access research and large language models for credible, trustworthy question answering**. In: *TPDL 2023*

CORE

# Reflections / limitations …

| ChatGPT | CORE-GPT |
|---|---|
| • Can get confused (esp. when answers are ambiguous) mixing content from entirely semantically different uses of a concept<br>• Can be made to argue your way producing biased text<br>• It can start inventing things / hallucinate … | • Answers need to be anchored to research papers.<br>• More honest about what it doesn't know => fewer hallucinations<br>• References make it easier to assess the trustworthiness of the answer. |

**Both**

- Powerful at synthesizing content and creating summaries
- Able to compare and contrast
- Can get confused (esp. when answers are ambiguous) mixing content from entirely semantically different areas / uses of a concept
- Can be made to argue your way producing biased text
- Critical thinking and judgement needs to be exercised

# CORE - AI Expert Finder

Prototype tool to automatically identify domain experts based on publications in >34m research papers

**Applications in:**

| Peer review | Proposal review | Consultant/Expert recruitment |

**Evaluation:**

- **Relevancy** - was the suggested candidate a suitable match?
- **Prior Knowledge** - was the suggested candidate previously known to the enquirer?
- **Conflict** - are there any conflicts of interest with the proposed candidate?

## Results

**74%** of suggested candidates were suitable

**34%** of suggested candidates were not known to enquirer

# CORE - AI Expert Finder

**Prototype tool to automatically identify domain experts based on publications in >34m research papers**

## Applications in:

| Peer review | Proposal review | Consultant/Expert recruitment |
|---|---|---|

## Evaluation:

- **Relevancy** - was the suggested candidate a suitable match?
- **Prior Knowledge** - was the suggested candidate previously known to the enquirer?
- **Conflict** - are there any conflicts of interest with the proposed candidate?

**Results**

**74%** of suggested candidates were suitable

**34%** of suggested candidates were not known to enquirer

# The crucial role of repositories in providing machine access to research content.

# Principle 1

Repositories should always establish a link from the metadata record to the item the metadata record describes using a dereferencable identifier pointing to the version held locally in the repository (if applicable). The dereferencable identifier should be provided in the appropriate metadata element in the used metadata format.

# Principle 2

Repositories should provide universal access to machines with the same level of access as humans have. It should be possible for machines to harvest the entire content of the repository in a reasonable time to enable a machine to maintain up-to-date information about the content held in the repository.

# Functional OAI-PMH endpoint

Use an external system to see how your repository is seen from the outside of your organisation.

**Test,** don't take that it works for granted

**Monitor:** the fact that it works now doesn't mean it can't go wrong when you least expect it

Overview

Harvesting

Content

OA compliance

DOI

Plugins

Membership

Settings

Start tutorial

**General information**

Last successful updating
28/01/2023

Total harvested outputs
55.25K

37%  20.26K
Full texts

Harvested with 27,323 issues affecting 36,347 records

**Harvesting issues**

ALL  ERRORS  WARNINGS  OTHER

⚠️ **Embargoed full text**

The full text download URL has restricted access. If the fulltext is intended to be embargoed or restricted in some way, no further action is required.

**8914** records are affected by this issue

💬 **Recomendation**

No action needed. However, you might use this to check if your embargo settings are valid.

DOWNLOAD IN CSV   SEE THE LIST

CORE

# Robots.txt

- Be careful not to block robots

- Don't give preferential treatment

```
# robots.txt for http:/          / …
# Indiscriminate automated downloads from
#   this site are not permitted
# See also: http:/        /RobotsBeware.html
# $Date: 2012/04/27 15:58:32 $
User-agent: *
...
Disallow: /pdf/
Disallow: /html/
...
User-agent: Googlebot
...
Allow: /pdf
Allow: /html
...
User-agent: Yahoo! Slurp
...
User-agent: msnbot
Crawl-delay: 20
...
Allow: /pdf
Allow: /html
...
```

CORE

# Validate metadata

- Adopt a relevant application profile (e.g. RIOXX.net)
- Use a metadata validation service, e.g. within the CORE Repository Dashboard

# Validate

- ⊞ Overview
- ↻ **Harvesting**
- 📄 Content
- ☑ OA compliance
- ▌▌▌ DOI
- ✽ Plugins
- ⚎ Membership
- ⚙ Settings
- ⏯ Start tutorial

## General information

Last successful updating
### 28/01/2023

Total harvested outputs
### 55.25K

**37%**
20.26K
Full texts

Harvested with **27,323** issues affecting 36,347 records

## Harvesting issues

| ALL | ERRORS | WARNINGS | OTHER |

⚠ **Embargoed full text**

The full text download URL has restricted access. If the fulltext is intended to be embargoed or restricted in some way, no further action is required.
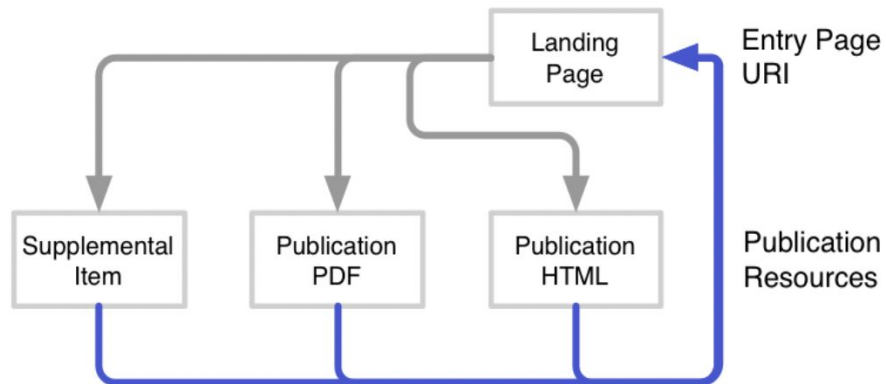
**8914** records are affected by this issue

💬 **Recomendation**

No action needed. However, you might use this to check if your embargo settings are valid.

DOWNLOAD IN CSV     SEE THE LIST

CORE

# Support Signposting

**Helping machines to navigate repositories in order to locate the content.**



CORE

# COAR Next Generation Repositories Working Group

Other · Open Access

## Next Generation Repositories: Behaviours and Technical Recommendations of the COAR Next Generation Repositories Working Group

Rodrigues, Eloy; Bollini, Andrea; Cabezas, Alberto; Castelli, Donatella; Carr, Les; Chan, Leslie; Humphrey, Chuck; Johnson, Rick; Knoth, Petr; Manghi, Paolo; Matizirofa, Lazarus; Perakakis, Pandelis; Schirrwagen, Jochen; Selematsela, Daisy; Shearer, Kathleen; Walk, Paul; Wilcox, David; Yamaji, Kazu

CORE

# Why is CORE important?

**Increase your contents' discoverability and prevent its misuse**
Search, Recommender, Discovery, PMC Linkout

**Make your papers uniquely identifiable and resolvable with PIDs**
OAI Resolver

**Assess and contribute to Open Access compliance and FAIRness**
Indexed by CORE badge

**Make your content machine readable**
Repository Health Check, CORE API, CORE Dataset, CORE FastSYnc,

**Become a CORE Member and benefit from lots more**
Dashboard: Metadata validation and monitoring

**>30M** monthly active users

# Next Generation Repositories: Behaviours

The next generation repository…

- manages and provides access to a wide diversity of resources, including published articles, pre-prints, datasets, working papers, images, software, and so on.
- is resource-centric, making resources the focus of its services and infrastructure
- is a networked repository. Cross-repository connections are established by introducing bi-directional links as a result of an interaction between resources in different repositories, or by overlay services that consume activity metadata exposed by repositories
- is machine-friendly, enabling the development of a wider range of global repository services, with less development effort
- is active and supports versioning, commenting, updating and linking across resources

CORE

# AI/ML for research intelligence and for improving repository workflows

# Affiliation extraction

## 1. Problem

Many metadata records do not have Some text …

Show an example how affiliations can be extracted. Show Grobid output …

How does this correspond with ROR
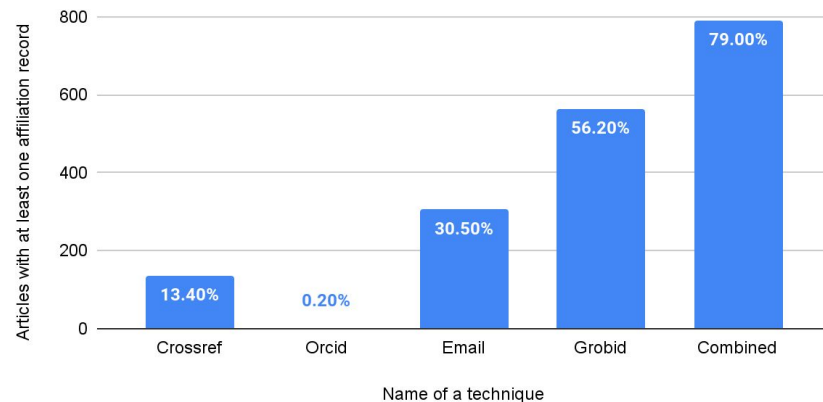
This is a problem we are currently working on

## 2. Publication footprint

CORE

# Affiliation extraction

- Many metadata records do not have affiliation data

- Affiliation is important for a range of use cases, including **publication footprint**

- At CORE, we developed a method to extract affiliation information from papers using a supervised ML model.

- Will propagate to the CORE API and Dashboard.

## Techniques comparison 1

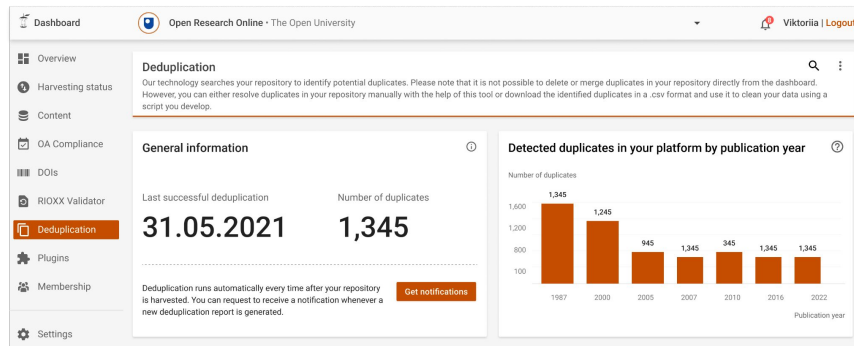Testing was performed on a sample of 1000 research papers in CORE



CORE

# Deduplication

How do duplicates look like and why do they occur in repositories?

| Example | Source Repository | Document Content | Why duplicates? |
|---------|-------------------|------------------|-----------------|
| A | Springer - Publisher Connector | Title = Profiling sugar metabolism during fruit ... | Exact same titles but documents aggregated from different repositories. |
| | ProdInra | Title = Profiling sugar metabolism during fruit ... | |
| B | Elsevier - Publisher Connector | Abstract = AbstractThe formation of smart, Metal Matrix Composite (MMC) structures through the use of solid-state ... | The abstracts are the same except for error introduced during document submission into different repositories. |
| | Loughborough University Institutional Repository | Abstract = This is an open access article under the CC BY license(http://\\ud\ncreativecommons.org/licenses/by/4.0/). The formation of smart, Metal Matrix Composite (MMC) structures through the use of solid-state... | |
| C | Swinburne Research Bank | Abstract = We present an analysis of ... 20-ms pulsars ... | Slight variation in text (20-ms vs 20 millisecond) on document versions on two different repositories. |
| | arXiv.org e-Print Archive | Abstract = We present an analysis of ... 20 millisecond pulsars ... | |
| D | Archivio della ricerca - Università degli studi di Napoli Federico II | Title = Simulation of Gaussian Processes and First Passage Time Densities Evaluation<br><br>Abstract= Motivated by a typical and .... first passage time probability densities. | Possibly different paraphrasing of the title for the exactly same abstract; the duplicates can only be identified when comparing "Abstract" rather than "Title". |
| | Archivio della ricerca - Università degli studi di Napoli Federico II | Title = Vectorized simulations of normal processes for first-crossing-time problems<br><br>Abstract = Motivated by a typical and ... first passage time probability densities. | |

CORE

# Deduplication

1. CORE uses an adapted version of locality sensitive hashing (simhash) for deduplication.
2. >90% F1-score performance
3. Deduplication powers our service including in the Dashboard for:
   a. versioning
   b. OA compliance (cross-repository)
   c. with affiliation extraction, this will allow us to warn institutions before outputs become non-compliant



Comparison mode

# Deduplication

### General information ⓘ

Last successful deduplication
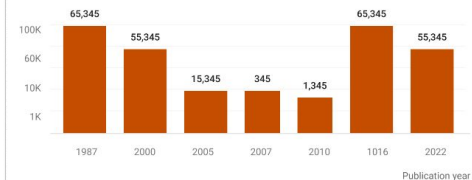
**31.05.2021**

Number of duplicates

**576**

Deduplication runs automatically every time after your repository is harvested. You can request to receive notifications whenever a new deduplication report is generated.

**Get notifications**

### Duplicates ⓗ

Number of duplicates

| Publication year | 1987 | 2000 | 2005 | 2007 | 2010 | 1016 | 2022 |
|---|---|---|---|---|---|---|---|
| | 65,345 | 55,345 | 15,345 | 345 | 1,345 | 65,345 | 55,345 |

**Comparison mode**

**List of possible duplicates**

| | | | | |
|---|---|---|---|---|
| | Zero and low carbon buildings: A driver for change in working practices and the use of computer modelling | 👁 LIVE IN CORE | Zero and low carbon buildings: A driver for change in working practices and the use of computer modelling | 👁 LIVE IN CORE |
| Repository | Open Research Online | | Open Research Online | |
| Author | Robina Hetherington, Robin Laney and Stephen Peake | | Robina Hetherington, Robin Laney and Stephen Peake | |
| DOI | 10.1109/iv.2010.86 | | 10.1109/iv.2010.86 | |
| OAI | oai:oro.open.ac.uk:21316 | | oai:oro.open.ac.uk:21316 | |
| Publication date | 21.09.2020 | | 21.09.2020 | |
| Deposited date | 30.10.202 | | 30.10.202 | |
| Version | Published | | Not available | |
| Abstract | This paper was selected for publication in MIT's Design Issues. The research takes an original approach by positioning experimentation as a comprehensive design methodology, rather than using the traditional... Show more. | | Not available | |
| Full text link | Unavailable | | Unavailable | |

| Duplicate | Different version | Not the same article |
|---|---|---|

← BACK

**COMPARE METADATA RECORDS**

| 2164/202 | Lorem ipsum dolor sit amet, consectetur adipiscing adipi | Lorem ipsum dolor sit | Need to be reviewed | 31/12/2019 | 👁 ⋮ |
|---|---|---|---|---|---|

ⓘ The below list contains the potential duplicates CORE identified. You can compare and review these potential duplicates and confirm them as duplicates or tell us that they are different. This will impact how CORE displays these articles in Search, API and other services. Specifically, by marking potential duplicates as different articles, these articles will be disassociated (they will not be part of the same **Work** entity).

**Possible duplicates in your repositories**

| OAI | Title | Author | Duplicate status ⓗ | Publication date ⓗ | | |
|---|---|---|---|---|---|---|
| 2164/202 | Lorem ipsum dolor sit amet, consectetur adipiscing adipis | Lorem ipsum dolor sit | Need to be reviewed | 31/12/2019 | 👁 | ⋮ |
| 2164/202 | Lorem ipsum dolor sit amet, consectetur adipiscing adipis | Lorem ipsum dolor sit | Duplicate | 31/12/2019 | 👁 | ⋮ |
| 2164/202 | Lorem ipsum dolor sit amet, consectetur adipiscing adipis | Lorem ipsum dolor sit | AM | 31/12/2019 | 👁 | ⋮ |

**DOWNLOAD CSV**

CORE

# Data enrichment

# Document classification

- Classification of research papers in a distributed environment is a problem.
- Established a benchmark for research document classification as part of the SDP/COLING conference.
- In the process of bringing themes to the CORE API.

# CORE moving to a membership model

August 2023

CORE will become an **independent** open scholarly infrastructure

CORE will **no longer** receive direct funding from Jisc

CORE will be **operated by** The Open University

**Membership**
(data providers)

**Sponsorship**

CORE

# CORE Membership

- A network of data providers who are committed to the ongoing success of the **Open Access movement**

- We provide **tools and benefits** for our members

- All CORE data providers are eligible to become CORE Starting Members **free** of charge
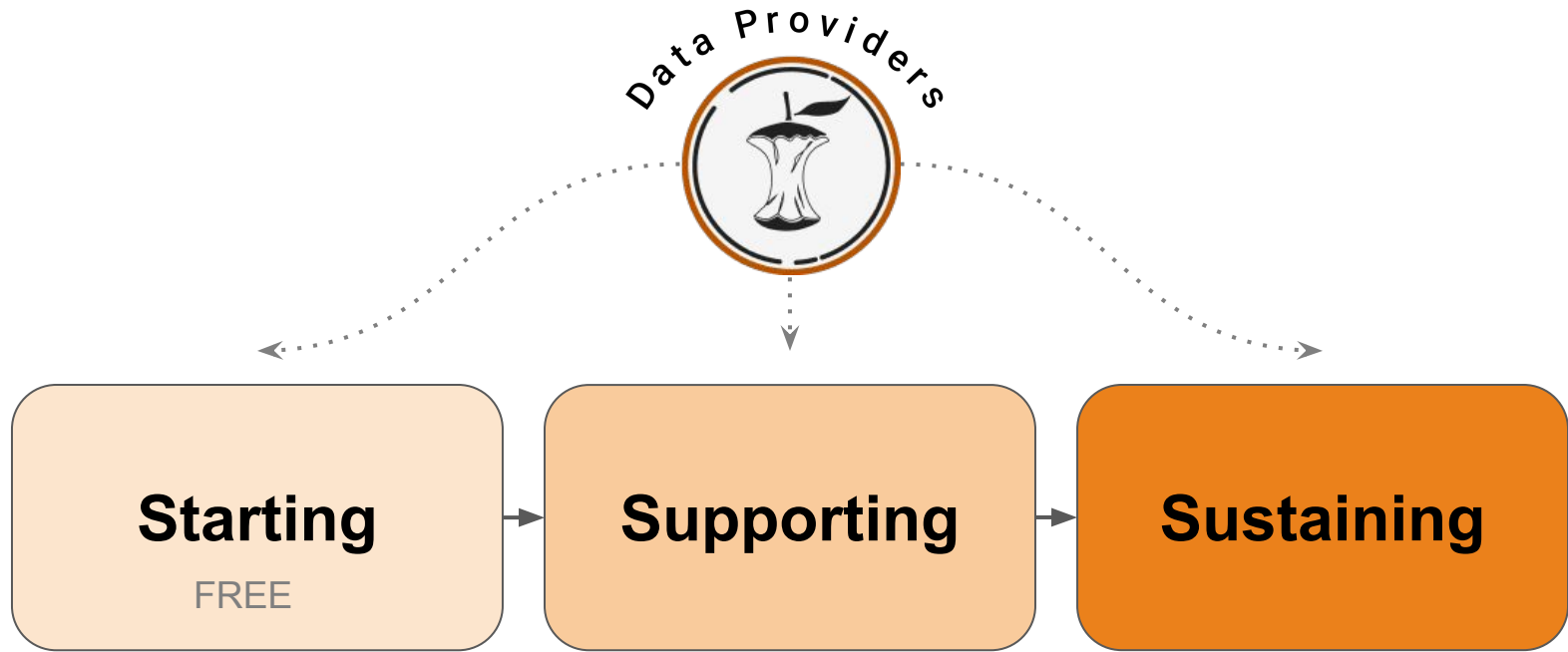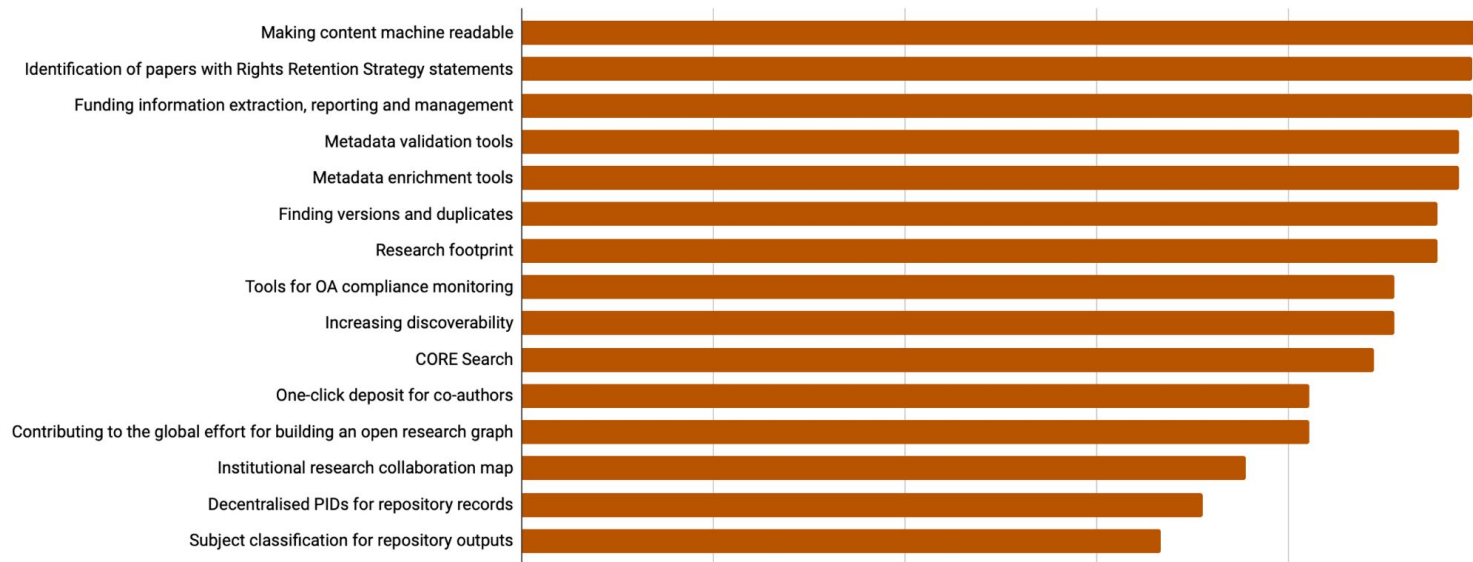
- Supporting and Sustaining Members:

    - help shape our development roadmap

    - support and sustain CORE

CORE

# Three levels of CORE Membership

# What matters to members (Board of Supporters survey)



| Category | |
|---|---|
| Making content machine readable | |
| Identification of papers with Rights Retention Strategy statements | |
| Funding information extraction, reporting and management | |
| Metadata validation tools | |
| Metadata enrichment tools | |
| Finding versions and duplicates | |
| Research footprint | |
| Tools for OA compliance monitoring | |
| Increasing discoverability | |
| CORE Search | |
| One-click deposit for co-authors | |
| Contributing to the global effort for building an open research graph | |
| Institutional research collaboration map | |
| Decentralised PIDs for repository records | |
| Subject classification for repository outputs | |

CORE

# More reading: references

Knoth, P. (2013). **From open access metadata to open access content: two principles for increased visibility of open access content**. In Open Repositories 2013. Retrieved from http://oro.open.ac.uk/37824/

Pride, D., & Knoth, P. (2020). **An Authoritative Approach to Citation Classification.** Proceedings of the ACM/IEEE Joint Conference on Digital Libraries in 2020. doi:10.1145/3383583.3398617

Kunnath, Suchetha N.; Pride, David; Gyawali, Bikash and Knoth, Petr (2020). **Overview of the 2020 WOSP 3C Citation Context Classification Task**. In: Proceedings of the 8th International Workshop on Mining Scientific Publications, Association for Computational Linguistics pp. 75–83.

Kunnath, Suchetha N.; Herrmannova, Drahomira; Pride, David; Knoth, Petr (2022). **A Meta-analysis of Semantic Classification of Citations** . Quantitative Science Studies, 2 (4), pp. 1170-1215

CORE

# More reading: references

Kusa, Wojciech; Hanbury, Allan; Knoth, Petr (2022). **Automation of Citation Screening for Systematic Literature Reviews using Neural Networks: A Replicability Study** . In: 44th European Conference on Information Retrieval, 10-14 Apr 2022, Stavanger, Norway Springer , 13185 , pp. 584-598

Nambanoor Kunnath, Suchetha; Pride, David; Knoth, Petr (2022). **Dynamic Context Extraction for Citation Classification**. In: The 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing, 20-23 Nov 2022, Virtual

Gyawali, Bikash; Anastasiou, Lucas; Knoth, Petr (2020). **Deduplication of Scholarly Documents using Locality Sensitive Hashing and Word Embeddings**. In: 12th Language Resources and Evaluation Conference, 11-16 May 2020, Marseille, France European Language Resources Association , pp. 894-903

CORE

# More reading: references

Óscar E. Mendoza, Wojciech Kusa, Alaa El-Ebshihy, Ronin Wu, David Pride, Petr Knoth, Drahomira Herrmannova, Florina Piroi, Gabriella Pasi, and Allan Hanbury. 2022. **Benchmark for Research Theme Classification of Scholarly Documents.** In Proceedings of the Third Workshop on Scholarly Document Processing, pages 253–262, Gyeongju, Republic of Korea. Association for Computational Linguistics.

Pride, David; Harag, Jozef; Knoth, Petr (2019). **ACT: An Annotation Platform for Citation Typing at Scale**. In: JCDL 2019 - ACM/IEEE-CS Joint Conference on Digital Libraries 2019, 2-6 Jun 2019, Urbana-Champaign, Illinois

Herrmannova, Drahomira; Pontika, Nancy; Knoth, Petr (2019). **Do Authors Deposit on Time? Tracking Open Access Policy Compliance** . In: 2019 ACM/IEEE Joint Conference on Digital Libraries, 2-6 Jun 2019, Urbana-Champaign, IL , pp. 206-216 BEST PAPER AWARD

CORE

# Take home …

- **ML/AI** has the potential to transform all stages of the research process, including how we carry out research, how we assess it and how we organise research knowledge.

- **OA** repositories play a key role in this process by providing machine access to research content.

- **AI/ML** already provides opportunities for improving the ways we use repositories, organise, enrich and curate content in them.

CORE

2023

THANK YOU