

| | |
|--------------------------------------|---|
| Acronyme | DISSEMIN |
| Titre de la proposition | Diffusion d'articles scientifiques et implémentation de politiques institutionnelles |
| Nom et prénom du coordinateur | Amarilli, Antoine |
| Organisme du coordinateur | Identification de l'établissement : Télécom ParisTech Adresse : 46 rue Barrault, 75013 Paris |
| Catégorie: | Web service development |
| Aide demandée | 99 360 euros |
| Coût complet : | 99 360 euros |
| Durée | 24 mois |
| Mots clés de la proposition | Open access, text and data mining, open repositories |

1. Description of the proposal

Enabling universal access to the scientific literature was the original reason why the World Wide Web was invented. Still, thirty years later, the problem is far from solved. University libraries are wrestling with academic publishers who ask for unaffordable subscription prices, and internet service providers block pirate websites which bypass paywalls.

The *open access movement* is an attempt to tackle the root of this issue by publishing articles openly on the web. To encourage this practice, universities and research funders have adopted policies that encourage their researchers to make their articles publicly available. Some editorial boards have moved away from traditional publishers and now run their journals electronically. Open access repositories have become the go-to platforms to read the latest research in many fields, even before it has been published in a peer-reviewed venue. Still, despite all these efforts, most of the scientific literature is still locked behind paywalls.

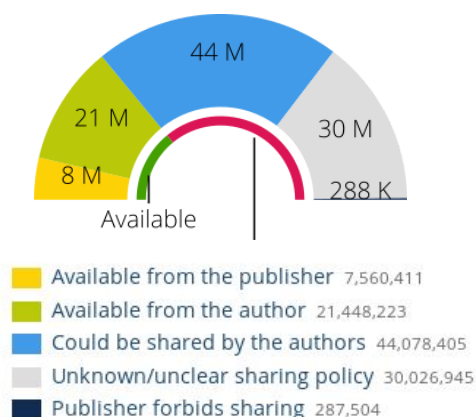
In this project, we propose to address this issue by developing the Dissemin platform. Dissemin (<https://dissem.in/>) is a web service to help researchers ensure that their publications are free to read. It does so by gathering information from various metadata sources and analyzes the online availability of each paper (over 100 million, see diagram) to determine which articles are freely available online and which are not.

Articles which are free to read can be either:

- openly available from the publisher itself ;
- or self-archived in an open access repository (known as green open access).

Other articles are broken down into three categories :

- those which could be legally uploaded by the authors to an open access repository, e.g., when the publisher's policy allows it;
- those for which the publisher's policy is unclear ;
- those which cannot legally be uploaded to an open repository.



In addition to this analysis, Dissemin also makes it easy for researchers to upload their publications to a selection of open repositories (HAL, Zenodo and soon OSF preprints). Depositing articles with Dissemin is simpler than doing it directly on the repositories, because Dissemin pre-fills all article metadata (authors, title, etc.) because it has already harvested it.

The Dissemin platform is already available online, but researchers currently have little incentive to use it: most researchers are not familiar with open access and do not realize the importance of making their research widely available online. In this proposal, our goal is to change this, using the policies adopted by research institutions and funding agencies: indeed, many such actors have adopted *open-access mandates* requiring their employees and project members to deposit their research openly, and compliance with these mandates is a prerequisite for promotions and grants.

Our goal is therefore to align Dissemin with institutional policies: we want to make Dissemin a tool of choice for universities and funders to implement their open access policies and monitor their uptake. This can be decomposed in several milestones that we detail below.

Objectives and milestones

> Integrating institutional repositories in Dissemin.

Institutions generally ask their researchers to use specific open access repositories to make their publications freely available, e.g., repositories hosted by the institutions, or national repositories such as HAL in France. It is therefore crucial that Dissemin interfaces with as many of these repositories to serve the needs of researchers. To do this, we need to develop connectors to various repository vendors, and to negotiate agreements with repositories to enable direct deposit interfaces, as well as maintaining and improving existing connectors.

> Linking publications to their funders and institutions.

There is currently no simple mechanism for funding agencies and research institutions to monitor the publications that are produced by the researchers that they employ or finance, and it is even harder to determine whether these publications have been deposited as open-access. For this reason, these actors cannot check the outcome of the open access policies that they adopt. We plan to simplify this analysis with Dissemin by recording the link between articles institutions, which will also encourage institutions to promote Dissemin to their users and to develop connectors for their own repositories. Establishing this link will be done using existing curated data sources (Fundref) and disambiguation of unstructured ones, using an entity linking algorithm. This entity linking algorithm constitutes by itself an important deliverable as it will be usable for other text mining applications. In particular, we intend to experiment with an email campaign in ESPCI to deposit the missing publications in HAL, following our existing collaborations with the ESPCI library (see “Consortium” below).

> Integrating funder policies and fostering their implementation.

Beyond linking publications to institutions, we want to represent the open access mandates adopted by these institutions, by integrating to Dissemin the existing databases for such policies (ROARMAP and SHERPA Juliet). Storing such policies in a structured format will allow us to assess their compatibility with publisher policies and with more general laws. For instance, the République




Numérique bill introduces gives researchers in French public research institutions the right to deposit their articles, overriding the usual embargo periods of many publishers.

State of the art

To our knowledge, no comparable platform to Dissemin currently exists. Tracking the publications authored within a given institution is generally done with Current Research Information Systems (CRIS), such as Pure (Elsevier) or Elements (Symplectic), none of which are open source and all of which can only be used internally in a single institution. These systems are therefore not suited to implement funder policies or represent overlapping affiliation layers which are common in research.

Existing achievements to justify the credibility of the proposal

Dissemin has already been used to run a campaign to deposit the publications cited in Wikipedia. After harvesting which research articles were cited in Wikipedia, the authors of these works were contacted and encouraged to use Dissemin to upload their publication to an open archive. Thanks to this campaign, thousands of Wikipedia references were improved by adding a link to an open access copy. The improvement of these references continues since this initiative, with a crowd-sourced effort run every year at the occasion of the international Open Access Week. There is therefore evidence that Dissemin is well suited for campaigns reaching out beyond single institutions, which would be impossible with existing CRIS systems.

| Repository: | |
|--|---|
|  | Zenodo is a general-purpose open repository hosted by CERN. If the document does not have a DOI yet, Zenodo will create one. <input checked="" type="radio"/> |
| Coming soon: | |
|  | Run by the Center for Open Science, OSF lets you share your preprints openly. You can also share supplemental data, materials or code associated with your article. <input type="radio"/> |
|  | HAL is a multidisciplinary open archive, sustained by the French Research Ministry. <input type="radio"/> |

We have already demonstrated our ability to integrate repositories into Dissemin via three approaches. First, we have already developed some connectors allowing Dissemin users to deposit their articles in a selection of open-access repositories directly from the Dissemin website. These repositories provide long term archival, are well indexed in search engines and are run by not-for-profit organizations.

On second form of repository integration is to provide data from Dissemin in institutional repositories. For instance, the University of Bordeaux uses Dissemin directly inside its institutional repository, using the Dissemin API which provides an open and machine-readable interface to Dissemin's data for third-party systems. Researchers of the university are presented with a list of publications that they could upload, directly from the institutional repository. This also simplifies the deposit process again thanks to the metadata exposed by Dissemin.

Third, integration with repositories can also take the form of self-hosted instances of Dissemin and/or contributions to the Dissemin codebase. In this spirit, we are working closely with a group of three German universities (TU Darmstadt, University of Stuttgart and TU Braunschweig) as part of their OpenIng project, funded by the German Federal Ministry of Education and Research (grant amount: 299k€). This project focuses on engineering sciences and intends to connect Dissemin to additional metadata sources, developing custom repository connectors and integrating with the authentication mechanisms of the partner institutions.

These three approaches illustrate how the project can realistically address its goals. They also underline the importance of the open source license chosen for the project: Affero GPLv3. This license ensures that third-parties can freely reuse and improve the platform, and that they must publish their own improvements to the platform under a similar license. This means that any improvement made to Dissemin in the scope of their project will be reusable on our own instance of the software. The successful funding of the OpenIng project, and the fruitful collaboration between their development team and ours, shows that this model offers a viable route to long term sustainability.

Overall, we believe that the Dissemin project can credibly address the goals of the call “ANR Flash sciences ouvertes”, by fostering the development of open science practices across many domain areas, and contributing to the development of open research data and metadata, in particular the open-access status, bibliographic metadata and institution/funder information for over 100 million research articles. We request financial help from ANR to fund the manpower required to accelerate the development of this platform and focus on new features: integrating institutional repositories, managing institution and funder information, and accounting for open-access mandates. We believe that the addition of these features to Dissemin can play an important role to develop open data and open science for research outputs.

Tasks

A. Data model design

To link publications and researchers to institutions and their open access policies, we have to design new data models to store this information in the platform. Our goals are threefold:

1. Leverage existing databases of institutions and policies, to learn from past experiences and make data ingestion easier.
2. Make it possible to infer the institutional policies applying to a given publication. Researchers can be affiliated, employed or funded by institutions in various ways, and policies have varying scopes and effects. The model should be able to represent these situations in a faithful way.
3. Minimize the curation load on the platform administrators and make it possible for users to easily fix and improve the data themselves. For instance, we want to make it easy for researchers to add their own research laboratory, or for librarians to improve the linkage of researchers to the institution they work for.

We plan to rely on the open knowledge base Wikidata, where a database of academic institutions is already being curated, drawing up on many established databases¹. Our goal is to mirror this dataset on the Dissemin platform while letting our users improve this data directly on Wikidata.

B. Institution linkage

Dissemin currently holds the metadata of more than 100 million publications from authors around the world. Our goal is to link as many of these publications to institutions, via various open data sources:

1. Curated databases such as Crossref’s funding registry (Fundref), linking 3.8 million publications to funders, or OpenAIRE’s funding relations. This data comes with institution identifiers that are linked to Wikidata.

¹ Specifically, for institutions: RNSR, GRID, ISNI, VIAF, SUDOC, AURÉHAL, RINGGOLD and others. For policies: ROARMAP and SHERPA/Juliet.

2. Free-text affiliation strings stored in Crossref, ISTEEX and other sources. These strings can be aligned to identifiers with entity linking tools, but possibly introducing noise in the data.
3. Employment information from ORCID, which come with time boundaries and institution identifiers but whose accuracy is also not perfect.

The goal is to perform this linkage without compromising the accuracy: coverage can then be optimized for a selection of institutions, when the interest develops.

C. Outreach and campaigning tools

Once we have sufficient publication coverage for a given institution, we want to make it easy for this institution to use Dissemin to improve the availability of these publications. To this end, we will run a pilot project in ESPCI to establish the appropriate workflows. This means:

1. Providing simple ways to export Dissemin's data, for instance by downloading a spreadsheet of all publications in a given institution, with their status and metadata;
2. Writing tutorials to help librarians and research administrators work with this data, relate it to their own research information system and conduct email campaigns with it;
3. Develop other tools as suggested by the feedback received in the pilot project.

We also plan to work with Wikimedia Italy to explore the lessons learnt by their own campaign run with Dissemin, which focused on freeing the access to articles cited in Wikipedia.

2. Consortium

This project is coordinated by Antoine Amarilli (Télécom ParisTech), and the second project member is Pierre Senellart (ENS, Télécom ParisTech).

Antoine Amarilli (Télécom ParisTech) obtained his PhD in 2016 from Télécom ParisTech. His research in data management focuses on managing uncertain data, data provenance, and on efficient query evaluation. He has published in prestigious venues in database theory (ICDT, PODS), theoretical computer science (ICALP, LICS) and AI (IJCAI, JAIR). He has served two years as a program committee member for ICDT. A vocal advocate of open access, Antoine is a founding member and board member of the CAPSH association that manages the development of Dissemin.

Pierre Senellart (ENS professor, Valda team leader) obtained his PhD in 2007 from Université Paris-Sud and his Habilitation à diriger les recherches in 2012 from Université Pierre-et-Marie-Curie. He has published numerous articles in internationally renowned conferences and journals (PODS, SIGMOD, VLDB Journal, Journal of the ACM, etc.) He has been a member of the program committee and participated in the organization of various international conferences and workshops (including PODS, WWW, VLDB, SIGMOD, ICDE). His research interests focus around practical and theoretical aspects of Web data management, including Web crawling and archiving, information extraction, document engineering, and intensional data management. He is a contributor to the development of Dissemin, and is also the scientific referent of both DI ENS and Inria Paris for scientific publications, and as such well-aware of the related legal and practical issues.

In addition to these two members, the project will also involve the existing development team of Dissemin. The following members intend to help the project (on a voluntary basis):

- **Marie Farge**, spokesperson.
 - Emeritus research director, LMD-IPSL, CNRS & École normale supérieure, PSL.
 - Marie has a deep understanding of the scholarly communications ecosystem, and is a renowned open science expert with extensive contacts in a variety of fields. She has served on many policy boards, including at the European Research Council.
- **Antonin Delpeuch**, core developer.
 - PhD student, University of Oxford.
 - Specializing in natural language processing, he brings in the expertise for the required disambiguation algorithms. He was nominated Open Access Champion by SPARC Europe for his work on Dissemin.
- **Lucas Verney**, research engineer.
 - PhD student, École normale supérieure, PSL.
 - With his extensive experience on developing open source projects and his expertise in bibliographic management systems, Lucas is best placed to enhance Dissemin's usability, reliability and feature portfolio.
- **Catherine Kounelis**, deposit process supervisor
 - Head Librarian, Library & Center for Historical Resources, ESPCI Paris, PSL
 - Catherine knows what it takes to fill an institutional repository and will play a pivotal role in making Dissemin fulfill the needs of institutions by running a trial campaign at ESPCI.
- **Charles Paperman**, database specialist
 - Maître de Conférences, University of Lille
 - With his expertise on database systems, Charles is improving Dissemin's storage back-end to make it efficient and extensible.
- **Patricia Mirabile**, policy specialist.
 - PhD student, Lettres Sorbonne Université (Paris IV).
 - Patricia pushed for ENS's open access mandate and now serves on the French Open Science Council (COSO). With her expertise on open science and reproducibility, she is able to position Dissemin to match researchers' needs.

We also have links with the following partners for the development of Dissemin:

- **OpenIng project** (TU Darmstadt, University of Stuttgart and TU Braunschweig) funded by the German Federal Ministry of Education and Research (grant amount: 299k€). This project focuses on the engineering sciences in Germany and intends to use Dissemin to foster open access in this field. This involves connecting Dissemin to additional metadata sources, developing custom repository connectors and integrating their authentication mechanisms.
- **Centre pour la Communication Scientifique Directe** (CCSD, UMS 3668). This partnership convened and funded by Couperin Consortium enabled us to add HAL as a repository to Dissemin. We still maintain this interface in coordination with HAL's development team.
- **University of Bordeaux** has integrated Dissemin in their institutional repository via our API.
- Other partners : we already have a working relationship with other stakeholders which run services that we rely on or interface with : ORCID, Bielefeld Academic Search Engine, Crossref, the Open Access Button, Unpaywall and others.

3. Impact of the proposal

The proposal intends to develop the following deliverables :

- A. **Fostering the archival of publications**, either directly through the Dissemin platform itself or by prompting deposits elsewhere, and more broadly by raising awareness about the issue in the research community.

This is the main goal of the platform, and the need for open archiving of scientific publications has been widely expressed across many academic communities. We can estimate success for this goal by counting how many deposits were made via the Dissemin platform, but we note that our goal is not that researchers deposit specifically via Dissemin. In particular, we do not especially intend to optimize long term user retention on Dissemin, for instance; we will also consider that this objective is achieved if third-party tools can use the Dissemin API and dataset to encourage deposits on a different platform.

- B. **Providing a new tool for research institutions and funders** to monitor and implement the policies they adopt, integrating with the information systems that they already use such as ORCID and their institutional repositories ;

This concentrates the new features developed as part of this project. By creating an institutional view on our data, Dissemin will become usable by research institutions to track the availability of their research output and conduct campaign to improve it. We plan to develop and document simple workflows for institutions to run their own email campaigns using our service. Our goal is to make these campaigns as accessible as possible, requiring little technical expertise. We will also ensure that they do not require institutions to share personal data (e.g., lists of employees) with us.

- C. **Delivering open data and software** that can be reused for other purposes and integrated in third-party systems freely.

As a byproduct of the Dissemin project, we are curating an open dataset that can be reused for other purposes. For instance, our users have created the [“disseminr” R package](https://docs.ropensci.org/disseminr/) to query our API from scripts in the R language.² This makes it possible to search for publications and retrieve their self-archiving policy and the repositories where they can be downloaded, for instance. Other tools produced as part of this project include a machine learning algorithm to predict the topic of a document submitted to HAL³, and an entity linking system for Wikidata, trained on academic affiliations⁴. We anticipate that the project will continue to produce similar deliverables in the future.

4. Communication of the results

The results of the project will be communicated by various means :

² <https://docs.ropensci.org/disseminr/>

³ <https://github.com/wetneb/haltopics>

⁴ <https://github.com/wetneb/opentapioca> , described in [arXiv:1904.09131](https://arxiv.org/abs/1904.09131)

- Presentations given at international events on scholarly communications (LIBER 2015, Open Repositories 2016, Open Science Conference 2016, submission to Open Science Fair 2019);
- Events organized around the project itself (hackathons, advocacy workshops);
- Publications about how we addressed the technical challenges of building the platform;
- Outreach to researchers and librarians via email campaigns and social media.

The Dissemin code is released under the GNU Affero General Public License v3.0 as recommended by DINSIC⁵ and we provide an open API which can be used freely without registration. Our data sources are all open, meaning that the platform can be run independently from our own instance. The source code is archived in Zenodo (DOI [10.5281/zenodo.2571984](https://doi.org/10.5281/zenodo.2571984)) and Software Heritage.

The platform is run by a non-profit association, the Committee for the Accessibility of Publications in Sciences and Humanities (CAPSH), which funds the technical infrastructure using donations and occasional financial support from our partners. This organization is run openly. The statutes and minutes of its meetings can be found at <https://association.dissem.in/>. CAPSH will continue to coordinate the development of Dissemin when the project is over: we intend to continue maintaining the platform even after the end of the project, as we have already been doing for several years.

5. Auto-evaluation processes

We propose to evaluate our project on the following criteria :

- Coverage and quality of links between articles and institutions. These links are ultimately used to link articles to the relevant open access policies adopted by the funders and employers and we therefore need to ensure that they are reliable.
- Accuracy of access statistics for articles. This is important to provide an accurate view of the accessibility of an institution's research output.
- Usage of our deposit interface for repositories. Experience has shown that it can rise sharply when running targeted campaigns.

6. Budget

The project's budget will be handled by Télécom ParisTech. Our main expense is a research engineer who will be employed by Télécom ParisTech to carry out all proposed improvements. We also include additional costs, namely, paying a computer for the engineer, paying for web server hosting to power the Dissemin service, and paying for trips to conferences to present Dissemin or our technical findings.

We emphasize that, while the budget only covers Antoine Amarilli, Pierre Senellart, and the engineer, we will also involve in the project the members of the current development team of Dissemin (see "Consortium" above), even though they are not administratively listed in the budget.

⁵ <https://www.numerique.gouv.fr/publications/politique-logiciel-libre/pratique/>

| Personnel | Yearly cost | Person.months | Total cost |
|--|-------------|---------------|------------|
| Antoine Amarilli (PI, 20%) | | 4.8 | |
| Pierre Senellart (20%) | | 4.8 | |
| Research engineer | 65600 EUR | 15 | 82000 EUR |
| Travel & other direct costs | Unit cost | Number | Total cost |
| Equipment (e.g., engineer computer) | 1500 EUR | 1 | 1500 EUR |
| Web hosting (e.g., server rental) | 1250 EUR | 4 | 5000 EUR |
| Travel costs (e.g., conferences) | 875 EUR | 4 | 3500 EUR |
| A. Total direct costs | | | 92000 EUR |
| B. Indirect costs (8% overhead): | | | 7360 EUR |
| Total estimated eligible costs (A+B): | | | 99360 EUR |
| Requested ANR contribution: | | | 99360 EUR |

7. Action plan

Timeframe of the project: January 2020 — December 2021

| Task | Members | Timeframe |
|-----------------------------|--|-------------------------------|
| 1 - Data model design | L. Verney, C. Paperman | Jan—May 2020 |
| 2 - Software development | Research engineer, L. Verney | Apr 2020 — July 2021 |
| 3 - Documentation | Research engineer, L. Verney | Jun 2020 — July 2021 |
| 4 - Assessment | L. Verney, A. Delpeuch, P. Mirabile | Sep—Oct 2020; Jul—Aug 2021 |
| 5 - Communication, outreach | M. Farge, A. Amarilli, C. Kounelis | Oct 2020 — Dec 2021 |
| 6 - Project management | A. Amarilli, A. Delpeuch, P. Senellart | Jan 2020 — Dec 2021 |

1 - Data model design

- Improve the current database schema to enhance scalability and maintainability;
- Design the data model to represent institutions and their policies;
- Create interface mocks for institution and policy integration in the platform;
- Linking with the Wikidata community to synchronize data models.

2 - Software development

- Implement the new features as defined by the goals of the project;
- Test and document the code as part of the development process, and maintain the codebase
- Respond to user feedback to improve the workflows.

3 - Communication, outreach

- Run focused email campaigns in partner institutions to encourage the deposit of articles;
- Present the project and the platform in open science conferences;
- Participate in policy discussions at the national and European level;
- Communicate with users and report the feedback to development team.

4 - Assessment

- Run evaluation campaigns to evaluate the quality and coverage of the curated data;
- Monitor and advise development workflows;
- Conduct focused user studies to identify usability issues;
- Communicate assessment results to the team for funding reports and external communication.

5 - Documentation

- Write documentation explaining the general architecture of the system;
- Write articles aimed at a wider audience about the technical challenges and solutions;
- If relevant, present the architecture in technical conferences.

6 - Project management

- Hire and supervise the research engineer;
- Report to funders and communicate with partners;
- Organize project meetings.

