

Les enjeux et la mise en place de liens entre les produits de recherche

Maud Medves (INRIA)
Alexis Arnaud (UGA)
Maria Grazia Santangelo (UGA)
Elias Chetouane (UGA)

Open Science Days - 2023-11-16



Plan

- 1 Contexte et enjeux des liens entre produits de la recherche
- 2 Initiatives existantes
- 3 Travaux en cours de développement
- 4 Synthèse

Plan

1 Contexte et enjeux des liens entre produits de la recherche

2 Initiatives existantes

3 Travaux en cours de développement

4 Synthèse

État des lieux

Increasing the positive impacts of journal policies on data sharing practices, Breakout 4, Research Data Alliance Plenary 20, Gothenburg, 2023

- 2016-2019 : augmentation du nombre de déclaration de disponibilité des données
- Modes principaux de diffusion des données :
 - ▶ 37% → 35% : dans un entrepôt de données
 - ▶ 23% → 15% : non indiqué
 - ▶ 10% → 16% : sur demande
- Increasing the positive impacts of journal policies on data sharing practices (video).
- Hrynaszkiewicz, I., Simons, N., Hussain, A., Grant, R. and Goudie, S., 2020. Developing a Research Data Policy Framework for All Journals and Publishers. Data Science Journal, 19(1), p.5. <http://doi.org/10.5334/dsj-2020-005>.

État des lieux

Data sharing practices and data Analysis availability upon request differ across scientific disciplines

- Étude de la disponibilité des “données critiques” (jeu de données, image, modèle) :
 - ▶ dans 875 articles de 9 disciplines
 - ▶ publiés dans Nature and Science
 - ▶ sur les périodes 2000-2009 et 2010-2019

- Disponibilité initiale :
 - ▶ complète 54.2% (33.0 – 82.8%)
 - ▶ partielle 71.8% (40.4 – 100%)
 - ▶ différences significatives : par domaine de recherche, type de données, journal, période de publication

- Tedersoo, L., Küngas, R., Oras, E. et al. Data sharing practices and data availability upon request differ across scientific disciplines. Sci Data 8, 192 (2021).
<https://doi.org/10.1038/s41597-021-00981-0>

État des lieux

Data sharing practices and data Analysis availability upon request differ across scientific disciplines

- 310 auteurs ont été contactés :
 - ▶ 39.4% (27.9 – 56.1%) de réponses positives
 - ▶ 19.4% de réponses négatives : fortes différences par domaine
 - ▶ délai médian de 15 jours pour obtenir les données
 - ▶ augmentation de la disponibilité globale des données après un délai de 60 jours :
 - ★ complète 69.5% (57.0 – 87.9%)
 - ★ partielle 83.2% (64.9 – 100%)
 - ★ plus faible disponibilité pour la période 2000-2009 (29.4% vs 56.0%)
- Lieux de diffusion des données :
 - ▶ 62.2% : supplementary materials
 - ▶ 35.0% : pas de lieu indiqué
 - ▶ 22.3% : archives de données
 - ▶ 19.7% : sur demande

 - ▶ certains jeux de données sont déposés dans plusieurs lieux

État des lieux

A study of the impact of data sharing on article citations using journal policies as a natural experiment

- 2 études sur l'impact de la diffusion des données sur le nombre de citations :
 - ▶ 17 journaux avec une politique de partage des données, 13 sans politique de partage
 - ▶ 200 articles (ou 2 années de publication) par journal
 - ▶ utilisation de Web of Science pour le suivi des citation sur les 5 ans après publication
 - ▶ pour la 2e étude : 2 journaux avec et 2 journaux sans politique forte de diffusion des données
- Résultats
 - ▶ étude 1 :
 - ★ pas de différence significative dans le nombre de citations entre avant et après la mise en place d'une politique de diffusion des données
 - ▶ étude 2 :
 - ★ 25 → 70% pour la diffusion des données entre avant et après le changement de politique
 - ★ 40% d'augmentation des citations après le passage à une politique forte
- Christensen G, Dafoe A, Miguel E, Moore DA, Rose AK (2019) A study of the impact of data sharing on article citations using journal policies as a natural experiment. PLOS ONE 14(12) : e0225883. <https://doi.org/10.1371/journal.pone.0225883>

Enjeux de lier les produits de recherche entre eux

- Pour soi :
 - ▶ Augmentation de la visibilité (et donc des citations) de ses publications.
 - ▶ Valorisation de tous les produits de recherche.
 - ★ **OSD@UGA 2022** : demi-journée consacrée à la prise en compte du logiciel dans l'évaluation des carrières en recherche.
- Pour les autres :
 - ▶ Nécessité pour la reproductibilité et la transparence.
 - ▶ Permet de découvrir des produits de recherche de proche en proche.
 - ▶ Valorisation de tous les membres d'une équipe de recherche.
- En prime : l'ANR et les financements européens imposent l'accès libre des publications et des données, de ce fait l'effort de diffusion étant fait, il est simple de faire des liens entre produits de recherche.

Plan

- 1 Contexte et enjeux des liens entre produits de la recherche
- 2 Initiatives existantes**
- 3 Travaux en cours de développement
- 4 Synthèse

Lier les publications, les données et les codes

Du côté des publications

Depuis HAL

Remplir les métadonnées “**Identifiants**” :

- codes : SWHID
- données : DOI

Identifiants

Identifiants

Ajoutez l'identifiant [DOI](#), [arXiv](#), [PubMed](#), [ADS](#), etc pour lier votre dépôt aux autres bases.

SWHID ▾



Données associées

Ajoutez l'identifiant [DOI](#) fourni par l'entrepôt où vos données sont archivées.



Lier les publications, les données et les codes

Du côté des données

Depuis un entrepôt - ex. : Recherche Data Gov

Remplir les métadonnées :

- **“Publications associée”** : HAL-ID ou DOI

Publication associée ?

Un ou plusieurs des champs suivants pourraient devenir requis si vous complétez l'un de ces champs optionnels.

Citation ?

Nom, Prénom (Année). Titre, Editeur. DOI



Type d'identifiant ?

Sélectionner...

Identifiant ?

ex. pour DOI : "10.15454/AEIOUY"

URL ?

Adresse URL, commençant par https://

Lier les publications, les données et les codes

Du côté des données

Depuis un entrepôt - ex. : Recherche Data Gouv

Remplir les métadonnées :

- **“Workflow de calcul”** : SWHID ou url vers Software Heritage

Métadonnées de workflow de calcul ^

Type de workflow ?

URL vers le dépôt de code externe ? +

Documentation ? +

Lier les publications, les données et les codes

Du côté des codes

Depuis Software Heritage

Remplir le fichier “**README**” :

- publications : HAL-ID ou DOI
- données : DOI



Software Heritage

Features

- Search
- Downloads
- Save code now
- Add forge now
- Help

README.md

-rw-r--r--

6.5 KB

codemeta.json

-rw-r--r--

3.2 KB

README.md

Siconos

ubuntu-latest release testing passing

release v4.4.0 License Apache 2.0

A software package for the modeling and simulation of nonsmooth dynamical systems in C++ and in Python.

Plan

- 1 Contexte et enjeux des liens entre produits de la recherche
- 2 Initiatives existantes
- 3 Travaux en cours de développement**
- 4 Synthèse

Projet FAIRCORE4EOSC

Cas concret de lien entre publications, données et codes

- Projet Horizon Europe financé sur la période juin 2022 à mai 2025
- En phase de spécification technique
- La phase de beta release démarrera fin novembre 2023

Projet FAIRCORE4EOSC

Présentation du projet

- 23 partenaires dont Inria via Software Heritage
- Objectifs : développer et mettre en oeuvre neuf composants de base de l'EOSC. Pour améliorer la découvrabilité et l'interopérabilité d'un nombre toujours croissant de produits de la recherche.
- Les composants de base : 3 exemples
 - ▶ MSCR : Metadata Schema and Crosswalk registry
 - ▶ SWHM : Miroir de l'archive Software Heritage
 - ▶ RSAC : Research Software APIs and Connectors

Projet FAIRCORE4EOSC

Work package 6 - *Research Software APIs and Connectors* (RSAC)

- Le WP6 et le travail mené par les RSAC
- Qui sont les RSAC ?
 - ▶ Deux archives ouvertes : InvenioRDM au CERN, Dataverse à DANS,
 - ▶ deux éditeurs : Dagstuhl et Episciences,
 - ▶ deux agrégateurs : swMATH et OpenAire.

Projet FAIRCORE4EOSC

RSAC et le rapport *Scholarly Infrastructures of Research Software (SIRS)*

- La base de travail : le rapport SIRS (scholarly infrastructures of research software) :
- Les quatre piliers : archiver le code source, le référencer, le décrire, le citer.
 - ▶ Archiver => s'assurer qu'on peut accéder au code source,
 - ▶ Référencer => s'assurer que l'on peut identifier les codes sources,
 - ▶ Décrire => faciliter la découverte du code source,
 - ▶ Citer => créditer les auteurs pour rendre le travail de développement attractif en matière de carrière.

Projet FAIRCORE4EOSC

Exemple de réalisation

- concrètement, ça signifie quoi ?
- ex. d'une archive ouverte (Zenodo)
 - ▶ dépôt dans Software Heritage des codes sources déposés dans Zenodo ;
 - ▶ obtention du SWHID associé ;
 - ▶ exposition du SWHID dans la notice maintenue par l'archive ouverte
- ex. d'un agrégateur (swMath)
 - ▶ extraire les références à des codes sources présentes dans des publications et déclencher l'archivage de ce code source connu de l'agrégateur mais absent de Software Heritage ;
 - ▶ exposer le SWHID correspondant dans la notice du code source maintenue par l'agrégateur

- Développement de nouvelles fonctionnalités d'Episciences dans le cadre du projet FAIRCORE4EOSC (actuellement disponible en test)
- Possibilité d'associer du code source ou des données au manuscrit que l'auteur soumet à une revue (via un identifiant, swhid ou DOI).

The screenshot displays the Episciences interface for a manuscript submission. At the top, the 'EPIsciences' logo is visible, along with a 'preprod' badge and the user 'Maud Medvet'. The page title is 'Epijinfo Sandbox journal'. The main content area shows the title 'A Generic Formalism for Encoding Stand-off annotations in TEI' with a 'Preprint' badge. Below the title, the authors are listed as 'Javier Pose ; Patrice Lopez ; Laurent Romary'. The abstract text describes a proposal for encoding stand-off annotations in the TEI standard. Metadata includes the source 'HAL:hal-01061548v1', section 'Test Maud', and submission date 'November 8, 2023'. Keywords are 'TEI,textual resources,annotations,stand-off annotations,[INFO.INFO-CL]Computer Science [cs]/Computation and Language [cs.CL]'. Action buttons include 'Delete article', 'Download this file', 'See the document's page on HAL', 'Update metadata', and 'Abandon publication process'. A 'Linked publications - datasets - softwares' section at the bottom offers options to 'Add dataset', 'Add software', and 'Add publication'.

- Le code / les données sont visibles par les relecteurs dans le formulaire de relecture.

The screenshot displays the 'Epijinfo Sandbox journal' interface. On the left is a navigation sidebar with links for Home, Folder, Simple text Content page, Dashboard, Submit a document, My Account, JOURNAL, www.episciences.fr, Users, Journal, Mail, Website, and Statistics. The main content area shows a document submission page for 'Curated Archiving of Research Software Artifacts: lessons learned from the French open archive (HAL)'. The authors listed are Roberto Di Cosmo, Morane Guerpeter, Bruno Marmol, Alain Morlet, Laurent Romary, and Jocifra Satobnka. The document is submitted on September 18, 2023. Below the text, there is a 'Software' section with a file explorer interface. The file explorer shows a directory named 'zdc0f46 /' containing files: 'config', 'example', 'src', 'tests', 'gignore', 'AUTHORS', 'CHANGES', 'LICENSE', 'Makefile', and 'README.mf'. At the bottom, there are sections for 'Linked publications - datasets - software' and 'Rating'.

Mettre en valeur données et codes associés sur la page de l'article

- Le code / les données sont affichés sur la page de lecture de l'article (et mentionnés dans le pdf)
- Un progrès important par rapport aux déclarations de disponibilité des données (DAS en anglais) actuelles.

The screenshot shows the EpiSciences website interface. At the top left is the EPI SCIENCES logo. The main header reads "Epjinfo Sandbox journal". Below this, the article title is "Curated Archiving of Research Software Artifacts: lessons learned from the French open archive (HAL)". The authors listed are Roberto Di Cosmo, Morane Gruempeter, Bruno Marmol, Alain Morbell, Laurent Romary, and Jozsef Szadovszki. The article text discusses the importance of preserving research results and the challenges of archiving research software. Below the text, there are buttons for "Download this file", "See the document's page on HAL", and "Update metadata".

The "Softwares" section shows a file browser interface for "Software Heritage". It displays a list of files with columns for "File", "Mode", and "Size".

File	Mode	Size
config		
example		
src		
tests		
glignore	-rw-r--r--	38 bytes
AUTHORS	-rw-r--r--	322 bytes
CHANGES	-rw-r--r--	1.8 KB
LICENSE	-rw-r--r--	25.8 KB
Makefile	-rw-r--r--	426 bytes
STAMP.mv	-rw-r--r--	1.1 KB

At the bottom, there are buttons for "Add dataset", "Add software", and "Add publication".

Plan

- 1 Contexte et enjeux des liens entre produits de la recherche
- 2 Initiatives existantes
- 3 Travaux en cours de développement
- 4 Synthèse**

- L'indication du lieu de diffusion des données dans les publications n'est pas assez présente.
- Ce lieu de diffusion n'est pas systématiquement pérenne.
- Le lien entre produits de recherche améliore la visibilité de tous les produits de recherche.
- Pour le moment, les liens sont à la charge des équipes de recherche, via HAL et les entrepôts de données.
- Mais une structuration forte en cours au niveau européen, notamment pour reporter la charge des liens au niveau des éditeurs ou diffuseurs.

Merci de votre attention !

Site web de la science ouverte à l'UGA.