# Software Heritage and IPOL, a fruitful collaboration towards reproducible research

**Miguel Colom**
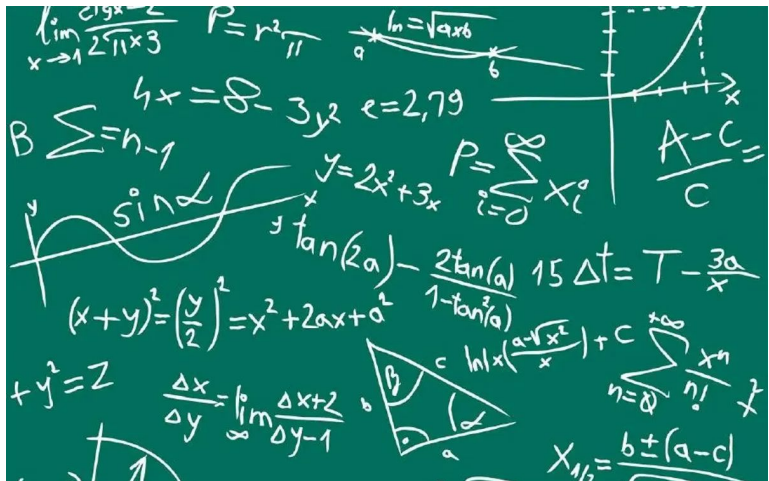miguel.colom-barco at ens-paris-saclay d0t fr

# The origin of the IPOL journal

- **Started** in October **2009**
- The initiative of the **Image Processing Group** at CMLA (now **Centre Borelli**) at ENS-Cachan (now **ENS Paris-Saclay**)
- **First article** published in **2010**

# IPOL's motivation

- The **reproducibility crisis** pointed out by Donoho et al
- Wanted to reveal the *real* **state of the art** in **image processing**
- Deeply understanding the methods. **All the mathematical details**
- **Reproducible research** :)

# How was it designed

- **Publication** = <span style="color:orange">**article**</span> + <span style="color:blue">**source code**</span> + <span style="color:green">**data**</span> as a whole
  - no longer a PDF + supplementary material (code, data)
- Let **everybody test** the **algorithms**
  - Scientific community, students, industry, …
  - With the **own data**
- Make it **straightforward to use it**:
  - **Online demos**

# Structure of a publication

- The **PDF** of the article
- The **peer-reviewed source code**, under a **FOSS license**
- Any associated **data**

- An **online demo** (supplementary material)
  - An **archive** of **experiments**

# Let's take a look

https://ipolcore.ipol.im/demo/clientApp/demo.html?id=201&key=ED5EFD91CC1FF7A8A9B4F7305C901F48

# Particularities in the editorial process (1 / 2)

- **Not easy** to find **reviewers** both **experts** in the **scientific field** and the **implementation details** (code)
    - **Our solution**: always **consider two reviewers**
    - **One** more **focused** on the **scientific aspects** of the article and the **other** on the **code**

    - → **Can we really separate those two aspects?**
        - **NO**. They need to **work together**. We can't simply split the tasks. Both the **article** and the **code** are part of the **same publication**.

# Particularities in the editorial process (2 / 2)

- The **editors work** with the **authors** to improve their **code** until it's published
- We need **permanent identifiers** and **pointers** to the code **during the review process**
- **Also after publication**:
  - The **sources** need to be **preserved**. **Permanent** storage
  - The **identifiers** needs to follow a **standard**. FAIR data
  - One should be able to **cite** the **whole** or **pieces** of the source **code**
  - The **sources** need to be **referenced**, with **different granularity levels**

# General difficulties related to Reproducible Research

- The **source code** in the **author's website** could **disappear**
  - For example, a researcher moves to another university
- The **project in Github** could be made **private**
- **Github** could **close**! (See the precedent of **Google Code**)
- The **author** could **alter the history** and the **commit's tree**
  - Several **tools available**: BFG Repo-Cleaner, git filter-repo, …
- What about a **DOI?** Same problem: the pointed object **can be altered**. **Integrity not ensured**. **Responsibility** on **publisher's side**
- Each forge might provide their **own non-standard formats** for **referencing** the code
  - Probably **not the adequate granularity**
- **Not** an **standard** way to **cite software**
  - HAL, IPOL, and others the **include** the **SWHID** though

# Quick note: *intrinsic identifiers*

- Some **identifiers** are "**extrinsic**":
  - **Not computed** from the **object itself**
  - For example: **the DOI**
- **Intrinsic identifiers** are also **based on the contents of the object**
  - For example: the **SHA-1** sum of a file

# So… what do we need?

- A **repository** of **all source code**, with **perpetual archiving**
- A **dynamic** archive
  - If a new commits arrive, we want then in the stored copy
- **Traceability** and **complete metadata**
- **Identifiers** at different **granularity** levels
  - **Intrinsic**
  - Be able to **cite the sources** in a **standard** way
    - In France: good solution with HAL + Software Heritage for citation of code
- **F**indable: **from the identifier** we should **arrive to the archive itself**
- **A**ccessible data: no registration, paywalls, …
- **I**nteroperable: an **open specification** of the identifiers
- **R**eusable: identifiers and formats we can **apply** in **other contexts**

# So… any *good* solution around?

- **Yes!**
- **Software Heritage** provides **all we need** to **evaluate** and **publish reproducible research** and conduct **open science**
  - **Permanent storage**
  - **Intrinsic identifiers** (**SWHID**)
    - Granularity: snapshot, release, revision, directory, file, line, …
  - **Open standard**: SWHID standardization in progress…
  - Possibility to properly incorporate it within **software citations**
  - **No cost** for **authors** or **institutions** to **use** the platform

## Regular crawling

These software origins get continuously discovered and archived using the listers implemented by Software Heritage.

Bitbucket
2,539,527 origins

56,983 origins

git
30,314 origins

R
26,984 origins

debian
136,866 origins

54,628 origins

GitHub
205,730,285 origins

gitiles
10,232 origins

GitLab
4,245,668 origins

+++ git
3,267 origins

Gogs
197 origins

GO
1,076,337 origins

Guix
50,149 origins

GNU
354 origins

heptapod
1,232 origins

launchpad
512,270 origins

Maven
312,428 origins

NixOS
48,590 origins

npm
3,595,535 origins

5,098 origins

Packagist
The PHP Package Repository
305,886 origins

fedora PAGURE
67,596 origins

Phabricator
201 origins

pub.dev
50,994 origins

python Package Index
524,009 origins

SOURCEFORGE
381,373 origins

stagit
318 origins

## Discontinued hosting

Discontinued hosting services. Those origins have been archived by Software Heritage.

GITORIOUS
122,014 origins

Google code {P}
790,026 origins

Bitbucket
336,795 origins

## On demand archival

These origins are directly pushed into the archive by trusted partners using the deposit service of Software Heritage.

eLife

HAL science ouverte

IPOL Journal

schema_version

object_id

`swh:1:cnt:41ddb23118f92d7218099a5e7a990cf58f1d07fa`

prefix

object_type

"snp" - snapshot

"rel" - release

"rev" - revision

"dir" - directory

"cnt" - content

origin_ctxt → `;origin=https://github.com/chrislgarry/Apollo-11`

visit_ctxt → `;visit=swh:1:snp:206c27c0c031c6aac6b5fedddba8fe082dea9836`

anchor_ctxt → `;anchor=swh:1:rev:3913f198f4383d4d638c0485d6aa902ff2f35828`

path_ctxt → `;path=/Luminary099/BURN_BABY_BURN--MASTER_IGNITION_ROUTINE.agc`

lines_ctxt → `;lines=64-72`

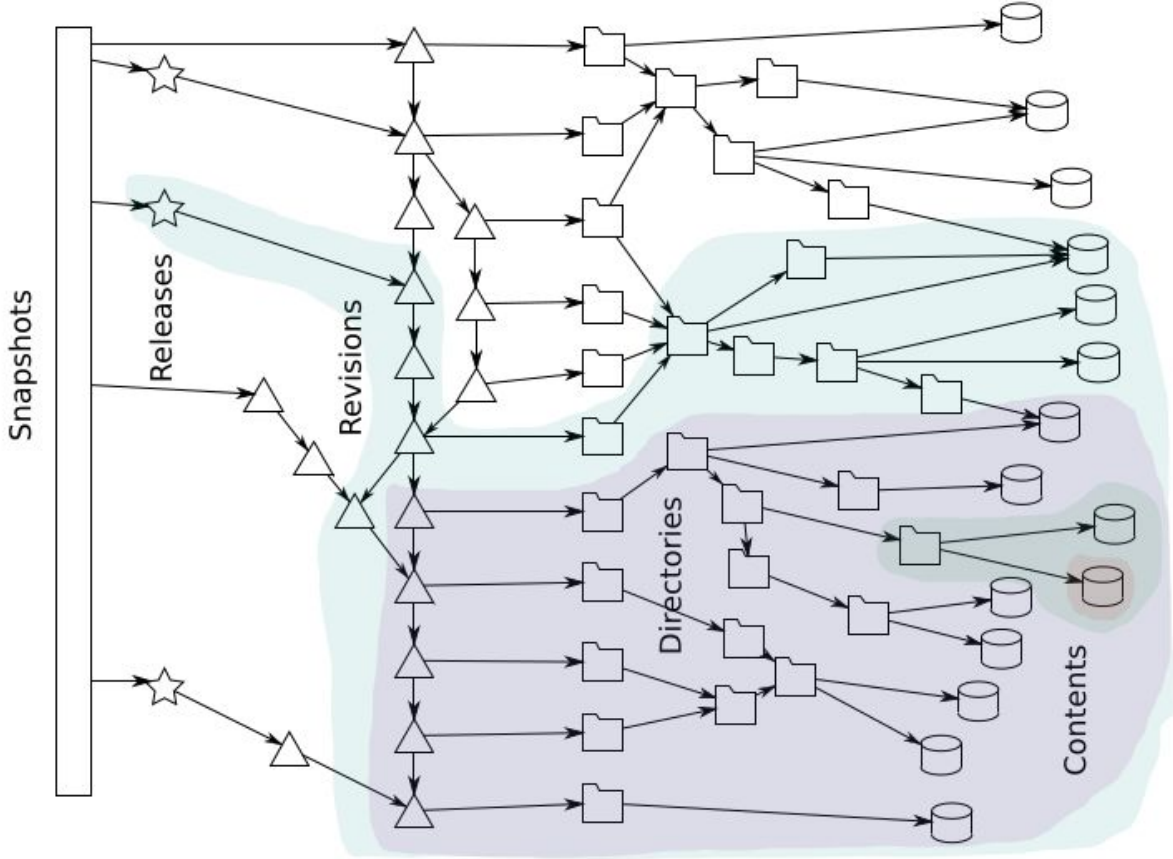# Current status

30+B SWHIDs in the Software Heritage archive

Mention in Linux Foundation's SPDX 2.2; IANA registered; WikiData P6138

# How does it work? Merkle tree

# IPOL Journal · Image Processing On Line

# Center/Surround Retinex: Analysis and Implementation

Jose-Luis Lisani, Ana-Belén Petro, Catalina Sbert

article  demo  archive

BibTeX info

*Communicated by* Jean-Michel Morel
*Demo edited by* Jose-Luis Lisani

## Abstract

The Retinex perception theory tries to mimic the human ability to cope with the high dynamic range of natural scenes. In 1986 E. Land proposed a formulation of this model in terms of a Center/Surround operation involving two steps, a local adaptation and a global transform. This model gave rise to the so-called Center/Surround tone-mapping algorithms. In this paper we unify the different Center/Surround algorithms proposed in the literature using a common framework and analyze several possibilities for the local and global operations involved.

## Download

- full text manuscript: PDF low-res. (576.4kB)  PDF (47.2MB) [?]
- source code: TAR/GZ  SWHID info </>

</> Software Heritage Archive

```
@softwareversion{sw-ipol.2021.391,
    title    = {{Center/Surround Retinex: Analysis and Implementation}},
    author   = {Jose-Luis Lisani, Catalina Sbert},
    date     = {2021-01-01},
    license  = {AGPL-3.0-or-later},
    version  = {1.6},
    swhid    =
{swh:1:dir:cd0313501fd340ef86219e2e52f6c8b202234d8e;origin=https://doi.org/10.5201/ipol.2021.391;vis
```

Copy to clipboard

# IPOL's workflow

1. The **author develops** and versions **with git** with Gitlab, Github, or any other collaborative platform
2. When the **code is submitted**, the **editors** take note of the **submitted revision** (commit ID)
3. **IPOL** might **create** a **public git repository** for the code if not available (the authors might submit a ZIP file, for example)
4. The **authors** can **continue developing**, but **IPOL freezes** at that **particular** submitted **revision**
5. In case of **changes** (typically bug fixes), the editors can **merge** after **reviewing** and **update** the version under review
6. When the publication is **accepted**, it's **submitted to Software Heritage** for **archival**

# IPOL's code publication

- **When the code is accepted, it's submitted to Software Heritage for archival**
  - At this moment: **manual process** by the copyeditor

  → **We're working on improving (automating) this…**

# IPOL's code publication: ideas for the short term

- **Automatic deploy** to **Software Heritage**. Not only after publication, but also during the review process
- **Use** of **SWHIDs** in the review process, whenever they're available
- **Automatic download** of the **sources from a revision** of the **git repository**. **No** more (controlled) **packages** from a particular revision
- By default **prefer** the **copy** of the sources **in Software Heritage** instead of the local copy, whenever it's possible
- Allow for **integrity checks**. For example, given a file we could compute its hash, compare to an IPOL's database of published codes, find it, and obtain its SWHID along with all the metadata. **Traceability**.

- In short, we need to use and **develop new tools** to ensure: **findability**, **availability**, **integrity**, **traceability**, and better **reproducibility** .

# Conclusion

- The **objective of IPOL** is **communicating reproducible research** on **algorithms**, with **detailed mathematical descriptions** and providing the **source code** under a **FOSS license**. **Open science**.
- The **inclusion** of **Software** as part of the **publication** is **not trivial**
  - More **complex review process**
  - Needs to **reference properly** the **sources** during the review process and after publication.
  - Needs **permanent archival**
- **Software Heritage** has proved to be an **excellent ally** for **IPOL**, since it provides a **complete solution** and **infrastructure**
  - This was expected: **software is not supplementary material**, but a **main research artifact**. Reproducible research needs to that **Software** is properly **referenced**, **archived**, and **cited**. **Software Heritage fully covered** this need.

# Thank you for your attention!