



Software Heritage

THE GREAT LIBRARY OF SOURCE CODE

Building sustainable infrastructures for (research) software

Roberto Di Cosmo

Founder and Director of Software Heritage
Computer Science Full Professor
Inria and Université Paris Cité

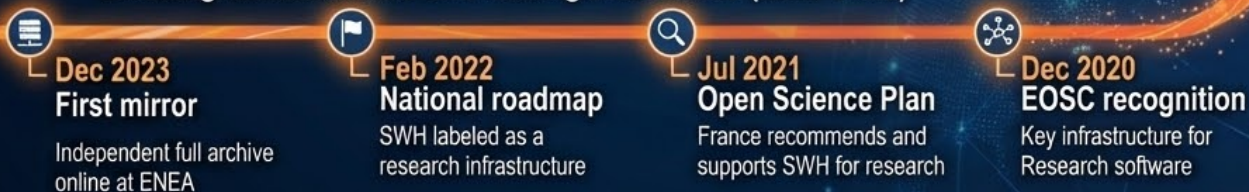


Software Heritage : a decade of commitment (2015-2025)

Laying the foundations (2015-2019)



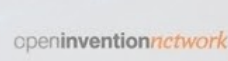
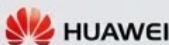
Building infrastructure and raising awareness (2019-2023)



Building solutions and strategic initiatives (2023-2025)



Inria





Software Heritage
THE GREAT LIBRARY OF SOURCE CODE

Inria with



The largest open source code archive

Unique digital common good built in France since 2015

Cultural Heritage

Industry

Research

Public Administration



Source files

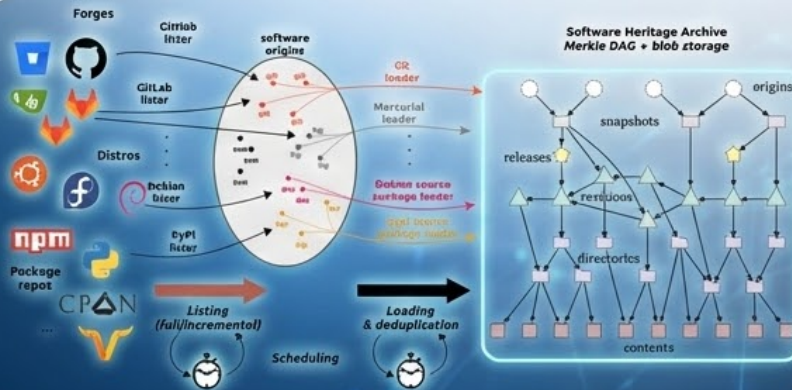
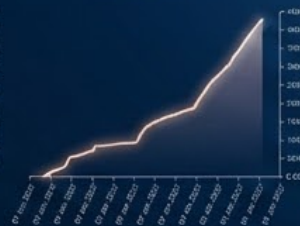
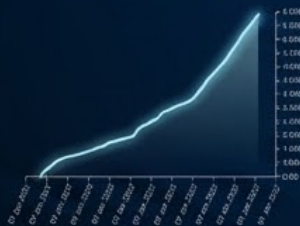
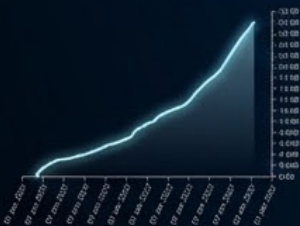
28,152,651,759

Commits

5,920,787,012

Projects

429,963,770



5000+ platforms

All versions, all history
development in a single graph

- 50 × 10⁹ nodes
- 800 × 10⁹ edges
~ 2 PB of storage

- ensures **availability**
- guarantees **integrity**
- allows **traceability**



of source
codes

**One infrastructure for multiple needs
replicated, open, multistakeholders**



Software Heritage
THE GREAT LIBRARY OF SOURCE CODE

The mirror network:

Creates a space for distributed ownership and collaboration



coherence
through
centralization



resilience
through
replication



Central Archive - Paris

UNIDUE - Duisburg-Essen

ENEA - Bologna

GRNET - Athens

- Central Node
- Active
- Deploying
- Agreement in Progress

Serving the scholarly ecosystem

Creating connections and mutualizing efforts



Software Heritage

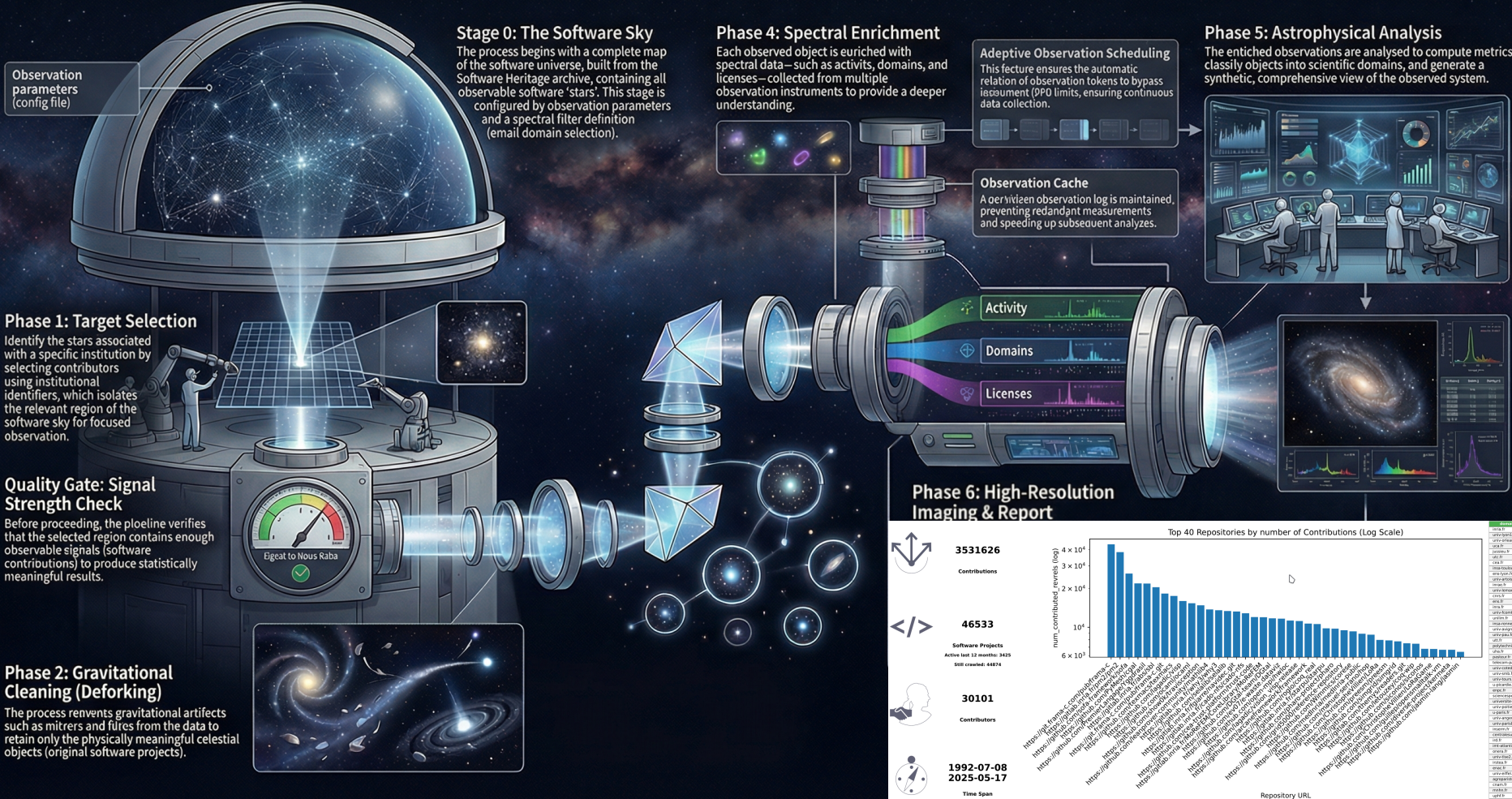
THE GREAT LIBRARY OF SOURCE CODE

Inria



unesco

From the Software Sky to Institutional Insight: Anatomy of an Observational Pipeline



Observation parameters (config file)

Stage 0: The Software Sky
The process begins with a complete map of the software universe, built from the Software Heritage archive, containing all observable software 'stars'. This stage is configured by observation parameters and a spectral filter definition (email domain selection).

Phase 1: Target Selection
Identify the stars associated with a specific institution by selecting contributors using institutional identifiers, which isolates the relevant region of the software sky for focused observation.

Quality Gate: Signal Strength Check
Before proceeding, the pipeline verifies that the selected region contains enough observable signals (software contributions) to produce statistically meaningful results.

Phase 2: Gravitational Cleaning (Deforking)
The process reverts gravitational artifacts such as mirrors and forks from the data to retain only the physically meaningful celestial objects (original software projects).

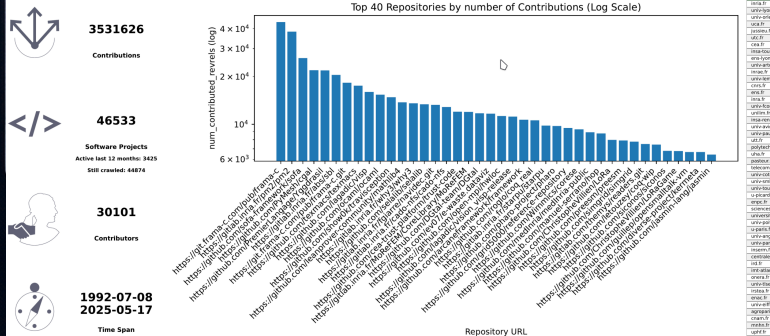
Phase 4: Spectral Enrichment
Each observed object is enriched with spectral data—such as activity, domains, and licenses—collected from multiple observation instruments to provide a deeper understanding.

Adeptive Observation Scheduling
This feature ensures the automatic relation of observation tokens to bypass incursion (DPA limits, ensuring continuous data collection).

Observation Cache
A persistent observation log is maintained, preventing redundant measurements and speeding up subsequent analyzes.

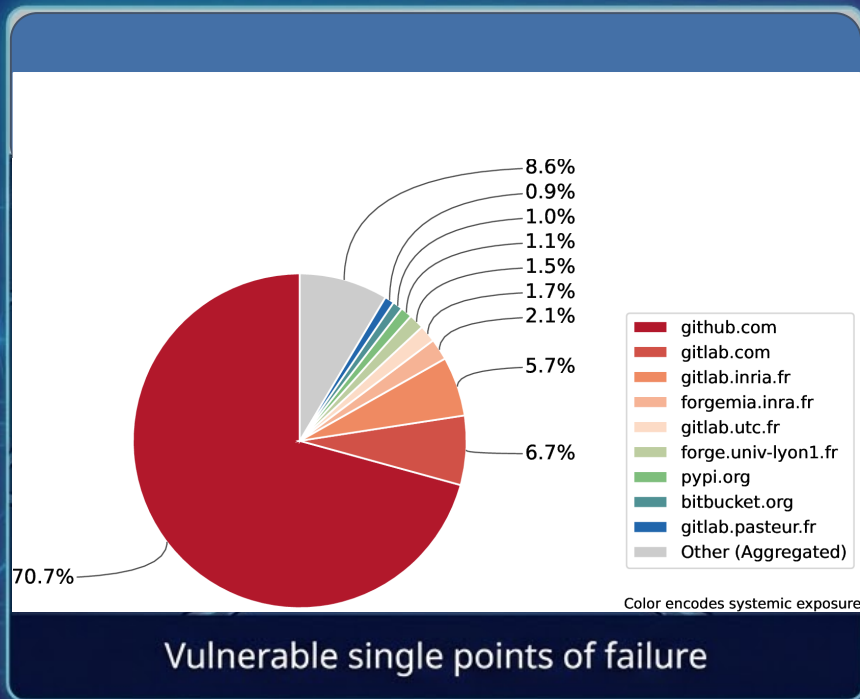
Phase 5: Astrophysical Analysis
The enticed observations are analysed to compute metrics, classify objects into scientific domains, and generate a synthetic, comprehensive view of the observed system.

Phase 6: High-Resolution Imaging & Report



The new frontier: ecumenical resilience

Global Concentration = Global Risk



The systemic response



Key principle: avoid balkanisation, mutualize core infrastructure