

Mining software repository with Pharo

03 July 2025

Nicolas Hlad

with Benoit Verhaeghe & Kilian Bauvent



ESUG2025 - Gdansk

Berger Levrault

 $(\mathbf{2}$

6

Introduction

What is Mining Software Repository ?

What are the specificities of Berger-Levrault regarding MSR ?

How we develop GitProjectHealth?

What can you do with it ? (demo)

What is next?



What is Mining Software Repository ?

(MSR)





Mining Software Repository **Definitions - 2**

- Today's collaborative development relies on Git Social Platforms (GSP) [Dabbish et al. 2012]
 - they are server implementations of Git with builtin social features
 - They contain valuable historical information over a software development





Fig - Commits distribution over time (GitHub)



Mining Software Repository
Definitions - 2

- Mining Software Repository (MSR) provides methods and tools to extract data from these platforms **[Hassan 2008]**.
- Among other, it allows us to:
 - Studying the impact of code smells [Steffen et al. 2010, Palomba et al. 2014]
 - Exploring developers code review [Bacchelli et al. 2013]
 - Predicting classes prone to change and defect [Bacchelli et al. 2010]
 - Retro-analysing a entire development process [Mockus et al. 2000]



Mining Software Repository
Existing tools

- General Mining data
 - PyDriller python tool for mining commit
 - Git-Miner Pharo tool, based on CLI-GitMiner
- Specific Mining Data
 - Javapers java lib mostly for Java file analysis in git repository (leveraging SPOON)
 - ModelMine retrieve UML model from project's artefacts
- Data Storage
 - Pandora —focus on long terme storage of Git social platform data
 - Software Heritage Archive of Git repository from GSP
- Computing Metrics
 - LinearB Productions and deployment metrics, data positioning with other companies



What are the specificities of Berger-Levrault regarding MSR



Industrial context

A quick word on Berger-levrault

- Berger-Levrault is
 - a group of international software editors
 - with divers sectors of activity (eduction, health, public administration, etc)
- The group acquired divers software editors over the past 30 years.
 - From different countries (France, Spain, Canada, Maroco, etc);
 - working with different technology (Java, C#, Typescript, Dart, etc);
 - and different Git Social Platform (Gitlab, Bitbucket, Azure Devops, etc).



Industrial context

A quick word on Berger-levrault - 2

- We use the project management system Altelissan's Jira to manage:
 - tickets (Bug, features, Hotfix, etc)
 - SPRINT (Agile development)
 - releases (delivering dates, testing software, etc).

Projets / Lab GL Dev Augmenté 総 …				
🕀 Résumé 😑 Chronologie 🛄 Tableau Kanban 🗠 Rapports 🖽 Liste	e 🗹 Tous les tickets 😑 Pages			
Q Recherche dans NH AH See O IB +7 Epic ~ Type ~ É	tiquette V Filtres rapides V		Regrouper: Re 2	
A FAIRE 24	EN COURS 22	TERMINÉ 2		
✓ Accélérer (35 tickets)				
Récupération des personnes par Equipe automatiquement MINING SOFTWARE REPOSITORY LABGL-156 ?	LABGL-91 Exploiter JaCoCo report Input les trucs modifié → output les tests associés (ou 0 si pas de test associé) ① 10 janv. 2025 ① LABGL-96	Add rule ensure unsubscribe AUTOCODEREVIEW FOR A WEMAGNUS 13 juin 2025 LABGL-223	••• ? 📧	

Fig - Kanban view of a SPRINT in Jira



Industrial context

A quick word on Berger-levrault - 3

• Our Jira and Git Social Platform environment are connected



Fig - Linking Jira Tickets to Commit and Merge activity in GitLab



How to mine from different Git Social Platforms (GSP) ? How to implement MSR metrics efficiently ? How to connect GSP data to Jira efficiently ?

We use Model Driven Engineering with Pharo-Moose



How we develop our solution with MDE

The conception of Git Project Health



Our MSR solution

GitProjectHealth

GitProjectHealth (GPH) is framework to extract and analyse data from Git Social platforms using Model-Driven Engineering (MDE).

tool for General Mining data & Metrics computing

- GitProjectHealth contributions are :
 - A unify model for all Git Social Platform
 - A framework to build custom metric from the model
 - A use of metamodel connector to extend any analysis to other platforms (e.g., Jira)



github.com/moosetechnology/GitProjectHealth



Main feature

GitProjectHealth

Key Features:

- Data Extraction & importation: Extracts data from major social platforms. Imports a model of specific Git entities.
- Visualization and Metrics: Visualizes data and computes metrics.
- Model Connection:

Connects models (e.g., Gitlab and Jira).





Project name	Project ID	closed merge request duration (second) - 22 June 2025	churn % in project (W=3) - 22 June 2025
Client	36189		19.49511162994309
inkscape	3472737	2505	1.5639652524938998





Git Model





Git Model





Git Model

2. New relations





Git Model

3. Concepts at fine granularity





API and Importer



Fig – APIs and GSP importers related to our Git model



Metamodel connection: Jira - Git Model



Metamodel connection: Jira - Git Model

🖸 GitProject-JiraConnector 🔺 💿 GitProjectJiraConnectorGenerator ! 🕨 instance sic 🔺 defineClasses			
GitProject-JiraConnector-Generator definition			
Filter v Filter overrides			
● All Packages ○ Scoped View ○ Projects ● Flat ○ Hier. ● Inst. side ○ Class side ● Methods ○ Vars <u>Class refs.</u> 🤍 <u>Implement</u>	ors 🔍		
© GitProjectJira(× ! Comment × ! *connectComr × * defineClasses × + Inst. side meth × < > 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	* >		
1 defineClasses			
2			
3 super defineClasses.			
4 ginmergekequest := self			
6 withPrefix: #GPH.			
<pre>7 glhCommit := self remoteEntity: #Commit withPrefix: #GPH.</pre>			
<pre>8 jiraIssue := self remoteEntity: #Issue withPrefix: #JP</pre>			
GitProject-JiraConnector Generator I instance sic A defineClasses			
GitProject-JiraConnector-Generator definition defineRelations			
Filter Filter			
● All Packages ○ Scoped View ○ Projects ● Flat ○ Hier. ● Inst. side ○ Class side ● Methods ○ Vars <u>Class refs.</u> Q <u>Implemente</u>	ors 🔍		
Pependencies × © GitProjectJira(× ! Comment × ! *connectComr × * defineRelation × < > 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	+ >		
1 defineRelations			
2	3 super defineRelations.		
2 3 super defineRelations.			
<pre>2 3 super defineRelations. 4 ((glhMergeRequest property: #jiraIssue) comment: 5</pre>			
<pre>2 3 super defineRelations. 4 ((glhMergeRequest property: #jiraIssue) comment: 5 'the jira issue associated to this merge request') 6</pre>			
<pre>2 3 super defineRelations. 4 ((glhMergeRequest property: #jiraIssue) comment: 5 'the jira issue associated to this merge request') 6 -<> ((jiraIssue property: #mergeRequest) comment: 7 'A merge request in Git associated to this issue')</pre>			
<pre>2 3 super defineRelations. 4 ((glhMergeRequest property: #jiraIssue) comment: 5 'the jira issue associated to this merge request') 6 -<> ((jiraIssue property: #mergeRequest) comment: 7 'A merge request in Git associated to this issue'). 8 ((glhCommit property: #jiraIssue) comment:</pre>			
<pre>2 3 super defineRelations. 4 ((glhMergeRequest property: #jiraIssue) comment: 5 'the jira issue associated to this merge request') 6 -<> ((jiraIssue property: #mergeRequest) comment: 7 'A merge request in Git associated to this issue'). 8 ((glhCommit property: #jiraIssue) comment: 9 'The jira issue associated with this commit')</pre>			
<pre>2 3 super defineRelations. 4 ((glhMergeRequest property: #jiraIssue) comment: 5 'the jira issue associated to this merge request') 6 -<> ((jiraIssue property: #mergeRequest) comment: 7 'A merge request in Git associated to this issue'). 8 ((glhCommit property: #jiraIssue) comment: 9 'The jira issue associated with this commit') 10 *-<> ((jiraIssue property: #commits) comment:</pre>			

Metamodel connection: Jira - Git Model







Quick demo



Our usage of *GitProjectHealth*

Deploying GPH at Berger-Levrault



Using GPH at Berger-Levrault Computing MSR Metrics

- We build a metric framework within GPH
 - They are either Projet or User centric
- For each Metric,
 - it loads entities from a time period (i.e., two dates)
 - it calculates the metric over a time windows (i.e. a Day, a Week, a Month, or a Year).
- 47 metrics are implemented so far.



Fig - Running all metrics in GitProjectHealth from a playground (simplified)



Fig - UML representation of the Metrics in GitProjectHealth

Using GPH at Berger-Levrault

Metrics computed every weeks (from 2024)



N° of commits should indicate whenever a developer is enough satisfy with a feature development progress

More commits **can** reflect more activity. However, developers can commits as much as they want, so this metric **can not be considered alone.**



Tendency

The number of contributor during the week. It represents the number of unique contributor to the project for the selected period.



Code Churn represent stability of introduced change Less is better Expected less than 7% for old project Less than 14% for new project



Close Merge Requests represent defective work. Lower the best. Closed MR merged MR, represent ended work. Created MR represent started work.

When they are close one to the other, it represents that we do not keep lot of unfinished work.

Merge Request metrics



Merge merged by dev - LinearB 2025 - Merge by dev by week **More is better**





Average number of comments/notes by MR

A high number means a lot of discussion for many which is unexpected 0/1 means no comments are made, interrogation for code review practice



Conclusion & Perspectives

and end.





Conclusion Recap











- Addressing limitations
 - The difficulty of maintaining a global metamodel by investigating the generation of GSP submetamodels from their OpenAPI
 - Discuss the purpose of the measures and consider which measures correlate with a healthy project.
 - Evaluating GPH against existing tools (PyDriller, Git-Miner, etc)
- Evolution
 - From GPH model to source code model navigating from repository to Famix model
 - Build usable knowledge maps by detecting parts of the repository that have become unknown to developers.



GitProjectHealth https://github.com/Evref-BL/Gitlab-Pharo-API Pharo Gitlab API https://github.com/Evref-BL/Gitlab-Pharo-API Pharo BitBucket API https://github.com/Evref-BL/Bitbucket-Pharo-API Pharo Jira API https://github.com/Evref-BL/Jira-Pharo-API

Example using GitProjectHealth:

Heatmaps https://github.com/Marpioux/Gitlab-HeatMap



[Steffen et al. 2010] Steffen M. Olbrich, Daniela Cruzes, and Dag I. K. Sjùberg. 2010. Are all code smells harmful? A study of God Classes and Brain Classes in the evolution of three open source systems. In 26th IEEE International Conference on Software Maintenance (ICSM 2010), September 12-18, 2010, Timisoara, Romania. 1-10.

[Palomba et al. 2014] Fabio Palomba, Gabriele Bavota, Massimiliano Di Penta, Rocco Oliveto, and Andrea De Lucia. 2014. Do They Really Smell Bad? A Study on Developers' Perception of Bad Code Smells. In Proc. of the 30th International Conference on Software Maintenance and Evolution. 101-110

[Bacchelli et al. 2013] Alberto Bacchelli and Christian Bird. 2013. Expectations, outcomes, and challenges of modern code review. In Proc. of the 35th International Conference on Software Engineering. 712-721

[Bacchelli et al. 2010] Alberto Bacchelli, Marco D'Ambros, and Michele Lanza. 2010. Are popular classes more defect prone?. In International Conference on Fundamental Approaches to Software Engineering. Springer, 59-73.

[Mockus et al. 2000] Audris Mockus, Roy T Fielding, and James Herbsleb. 2000. A case study of open source software development: the Apache server. In Proc. of the 22nd international conference on Software engineering. Acm, 263ś272. [Hassan 2008] A. E. Hassan, "The road ahead for Mining Software Repositories," 2008 Frontiers of Software Maintenance, Beijing, China, 2008, pp. 48-57, doi: 10.1109/FOSM.2008.4659248.

[Dabbish et al. 2012] Dabbish, L., Stuart, C., Tsay, J., & Herbsleb, J. (2012, February). Social coding in GitHub: transparency and collaboration in an open software repository. In *Proceedings of the ACM 2012 conference on computer supported cooperative work* (pp. 1277-1286).



Annexe



Titre de section

Texte du titre

Open-Source Analysis (github)

- Analyzed 30 days of activities from Eclipse, Microsoft, and MooseTechnology.
- Analyzed ~4457 entities over 30 days
- Visualizations: daily commit distributions, user activity



A user's commit activity by day, during the month of September 2024, on **moosetechnology**

Туре	Occurrences	Occurrences %
Anomaly	19	27%
Task	21	30%
Epic	26	37%
None	4	6%
Total	70	100%

Issues occurrences in September 2024 for **WeHR**

Industrial Case at Berger-Levrault (gitlab)

Connected Jira model with Git model to analyze merge request distributions across different issue types.

• 27% of Merge Requests linked to bug-related Jira issues.



Commit activity on vscode - around September 2024





Sunburst: last activity on the code base



Fig — Sunburst representation of a developer activity in a project