

# REPertoire: A Cross-System Porting Analysis Tool for Forked Software Projects

Baishakhi Ray, Christopher Wiley, Miryung Kim  
The University of Texas at Austin  
{rayb, thewiley, miryung@ece}.utexas.edu

## ABSTRACT

To create a new variant of an existing project, developers often copy an existing codebase and modify it. This process is called software forking. After forking software, developers often port new features or bug fixes from peer projects. REPERtoire analyzes repeated work of cross-system porting among forked projects. It takes the version histories as input and identifies ported edits by comparing the content of individual patches. It also shows users the extent of ported edits, where and when the ported edits occurred, which developers ported code from peer projects, and how long it takes for patches to be ported.

## Categories and Subject Descriptors

D.2.7 [Software Engineering]: Distribution, Maintenance, and Enhancement

## Keywords

software evolution, forking, porting, repetitive changes, code clones

## 1. INTRODUCTION

Software forking occurs when a developer or a group of developers splits off software into separate conceptual entities by copying and modifying an existing project. Forking is particularly common in free and open source software projects, where differing visions and personality clashes occur without an unifying profit motive. For instance, the split of FreeBSD and NetBSD from 386BSD, XEmacs from GNU Emacs, and LibreOffice from OpenOffice are well known forks. Proprietary software may also be forked to support the needs of multiple customers with different feature requirements.

Forking is often considered to be counter-productive. As multiple peer projects evolve in parallel, developers need to take similar features or bug fixes from one project to another [2]. Such cross-system porting practice often incurs

duplicate maintenance effort. This paper presents REPERtoire, a tool that analyzes the extent and characteristics of cross-system porting. It allows users to analyze the number of lines of code ported from the patches of peer projects, the developers responsible for those ported edits, the time taken to port patches from peer projects, etc. It presents the temporal and spatial characteristics of cross-system porting using various graphical views. It also supports interactive browsing of ported edits. Currently it is fully integrated with the state of the art version control systems such as Git and Mercurial. These analyses are designed to aid managers and architects to make informed decisions about the maintenance of forked software systems.

## 2. REPERtoire FEATURES

Suppose Sheryl is a manager working for the Exemplar corporation, which writes and sells software to enterprise customers. Two years ago, a particularly large customer requested a feature that required extensive modifications to the main product. To accommodate this customer's needs, the company forked the main product and made the necessary custom changes. Since then, a considerable amount of engineering effort has been continually spent to port bug fixes and security patches from the main product. Sheryl is contemplating whether it would be worthwhile to merge the two products back instead of duplicating maintenance effort.

Sheryl may need to analyze how the products evolve in parallel and how often cross-system porting occurs. She needs to know where the porting effort is focused on, who are the main developers porting code from peer projects, and how often cross-system porting happens, etc. She needs to know which directories and files mostly consist of ported edits. She may also be interested in knowing how long it takes for bug fixes and security patches to propagate from the main product to the other. These are the questions that REPERtoire can help Sheryl to answer. For presentation purposes, we refer to the main project and the forked project as *A* and *B* respectively in the following subsections.

**Porting Frequency View.** Suppose that Sheryl wants to know how often cross-system porting occurs. Given the version histories of *A* and *B*, REPERtoire visualizes the extent of code ported from one project to another over time. In the Porting Frequency View in Figure 1, the x-axis shows time in months and the y-axis shows the average percentage of ported edits with respect to total edits in each commit. Sheryl may select to see only the edits ported from *A* to *B*, only the edits ported from *B* to *A*, or both. Sheryl

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SIGSOFT'12/FSE-20, November 11–16, 2012, Cary, North Carolina, USA.  
Copyright 2012 ACM 978-1-4503-1614-9/12/11 ...\$15.00.

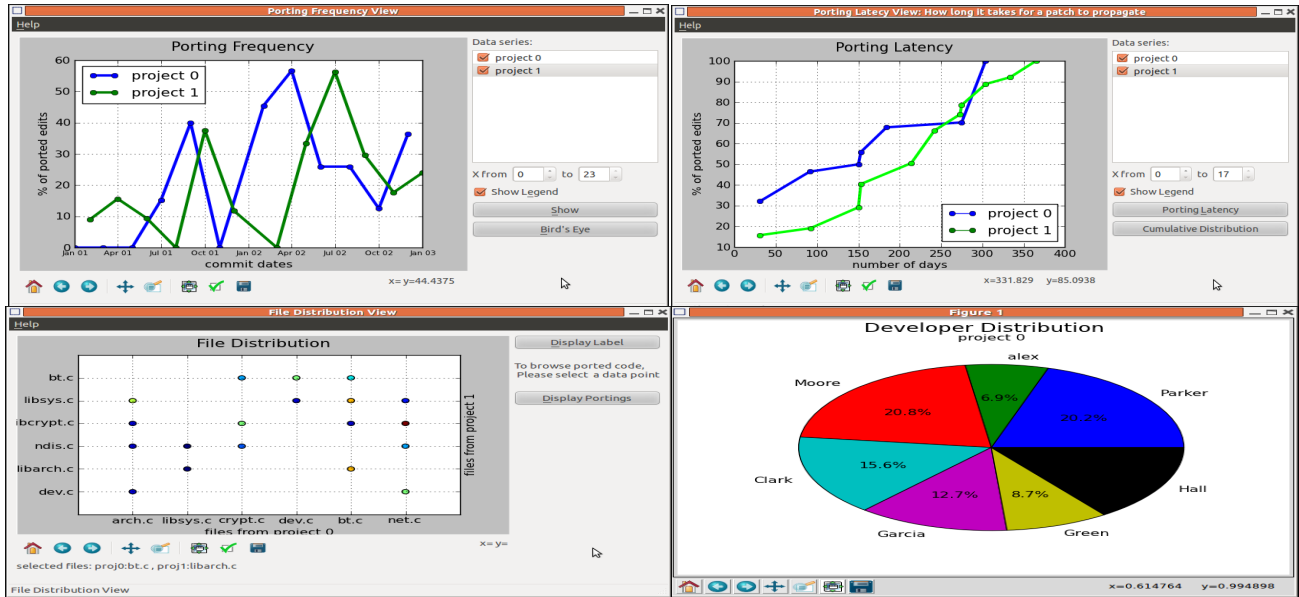


Figure 1: REPERTOIRE analysis of cross-system porting between two forked projects.

may see that 90% of the commits to  $B$  are ported from the patches of  $A$ , whereas 95% of the commits to  $A$  are not ported from  $B$ , indicating that most engineers working on  $B$  spend their time porting code and little time writing original code. On the other hand, if Sheryl notices that most edits to either system are not ported, then she may conclude that the systems are diverging further apart.

**File Distribution View.** To figure out where her organization is spending time porting code, Sheryl needs to see which pairs of files share ported edits between  $A$  and  $B$ . REPERTOIRE helps Sheryl by presenting the File Distribution View of the source and target of ported edits. This view is a scatter plot with files from  $A$  making up the x-axis and files from  $B$  making up the y-axis. A point is plotted at  $(x, y)$  if there is a ported edit from file  $x$  to file  $y$  or vice versa. Users can also grasp the amount of ported edits by inspecting the color of the dot. The darker the color is, the higher the ratio of ported edits to the total lines of code in the file. This allows the user to answer which files have the most ported edits and which files have the highest ratio of ported edits. The File Distribution View in Figure 1 shows an example of this file distribution view. By selecting any point on this view, Sheryl can browse all ported lines between the two files and investigate who ported the corresponding code, the commit dates, etc.

**Developer View.** To ask her team about the feasibility of merging  $A$  and  $B$ , Sheryl may need to identify the developers who have a deep understanding of both projects. A reasonable heuristic for finding such developers is to simply identify the developers who do a lot of porting work. REPERTOIRE displays a pie chart showing which developers are responsible for what fraction of the total ported lines. See the Developer Distribution View of Figure 1.

**Porting Latency View.** REPERTOIRE shows a user how long it takes for individual patches to propagate from one system to another system by presenting a cumulative distribution of porting latencies. The Porting Latency View

of Figure 1 shows the number of days between when a patch is first committed to one system and when a similar patch is committed to a target system.

### 3. IMPLEMENTATION & EVALUATION

REPERTOIRE analyzes *diff*-based program patches of two forked projects to identify the ported edits. It works in two phases. In the first phase, REPERTOIRE uses CCFinderX [3] to identify similar edit contents (clones) in the input patches. In the second phase, REPERTOIRE determines if two identified clones represent similar edit operations by comparing the edit operation types (i.e., addition, deletion, and modification) using an  $N$ -gram matching algorithm [1]. By comparing the commit dates of similarly edited code regions, REPERTOIRE disambiguates the source vs. target of the ported edit.

This tool demo paper expands on a tool that we developed to study co-evolution of BSD products and a detailed description of REPERTOIRE is described in [4]. The input wizard of REPERTOIRE gathers information about the version histories of forked projects. The analysis wizard of REPERTOIRE then visualizes cross-system porting analysis results between the input projects using several views: Porting Frequency View, Developer View, Porting Latency View, and File Distribution View. Using the inputs specified by the user in the input wizard, REPERTOIRE’s back-end extracts individual *diff*-based patches, developers, and commit dates from the version control repositories and compares the content and edit operations of the patches using CCFinderX.

Table 1 shows example inputs. After identifying cross-system ported edits, the back-end stores the results in a database, which can then be loaded from GUI visualization components. The internal structure of REPERTOIRE is shown in Figure 2.

According to our previous study [4], REPERTOIRE detects ported edits with 94% precision and 84% recall, at a token threshold 40. This accuracy was determined by comparing

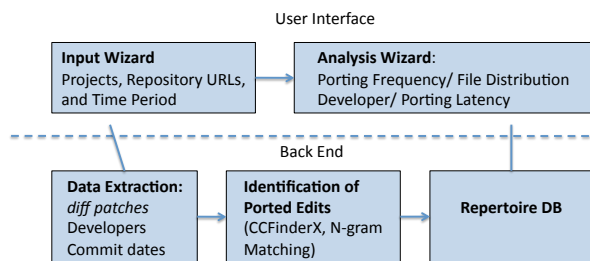


Figure 2: REPERTOIRE internal components.

	Input Types	Example Inputs
	working directory	<code>/var/tmp</code>
	CCFinderX path	<code>/usr/bin/ccfx</code>
Project 1	version control type	Git
	repository root	<code>/path/to/myrepo</code>
	time period	11/2/2010 - 9/31/2011
Project 2	version control type	Mercurial
	repository root	<code>/path/to/anotherrepo</code>
	time period	11/2/2010 - 9/31/2011

Table 1: Example inputs to REPERTOIRE.

REPERTOIRE’s result with a manually constructed ground truth set of ported edits on a sampled evolution period of OpenBSD releases 4.4 to 4.5.

In 18 years of parallel evolution of the BSD family, on average, FreeBSD ports 13.77% of edited lines from NetBSD and OpenBSD, while 15.52% and 10.74% of edited lines in NetBSD and OpenBSD originate from the other two BSDs respectively. 26.12%, 58.85%, and 44.85% of active developers in FreeBSD, NetBSD, and OpenBSD port patches from the other BSDs. The average time taken to port patches from peer projects in FreeBSD, NetBSD, and OpenBSD are 734, 725, and 944 days respectively. In all three projects, porting is mostly localized within 20% of the modified files.

## 4. SUMMARY

This paper presents REPERTOIRE that analyzes the extent of cross-system porting among projects forked from a common ancestor. Using REPERTOIRE, managers and engineers can measure the frequency of cross-system porting, learn which developers do how much of the porting work, investigate the trend of cross-system porting work over time, and the spatial distribution of ported edits with respect to the file system structure.

## 5. REFERENCES

- [1] G. W. Adamson and J. Boreham. The use of an association measure based on character structure to identify semantically related pairs of words and document titles. *Information Storage and Retrieval*, 10(7-8):253–260, 1974.
- [2] D. M. German, M. Di Penta, Y.-G. Gueheneuc, and G. Antoniol. Code siblings: Technical and legal implications of copying code between applications. In *MSR ’09: Proceedings of the 6th IEEE International Working Conference on Mining Software Repositories*, pages 81–90, 2009.
- [3] T. Kamiya, S. Kusumoto, and K. Inoue. CCFinder: A multilingualistic token-based code clone detection system

for large scale source code. *IEEE Transactions on Software Engineering*, 28(7):654–670, 2002.

- [4] B. Ray and M. Kim. A case study of cross-system porting in forked projects. In *FSE-20: ACM SIGSOFT the 20th International Symposium on the Foundations of Software Engineering*, 2012, to appear.

## APPENDIX

REPERTOIRE is an open source tool and can be downloaded from <http://dolphin.ece.utexas.edu/Repertoire.html>. This section describes the steps required to run REPERTOIRE.

### A. INSTALLATION

1. Install required libraries
  - Python 2.7
  - Qt 4.x: a cross-platform application and UI framework
  - pyuic4: a UI compiler for Qt that comes with the PyQt package.
2. Run ‘make’ from `src/`
3. Run ‘make’ from `src/analysis/`
4. Obtain a working copy of CCFinderX for intended platform
  - Ensure the execution of CCFinderX by running a command ‘`ccfx d cpp somefile.cpp`’

### B. POPULATING A DATABASE

REPERTOIRE takes the repository location and time period of version histories as input and identifies ported edits among the input projects. It requires a working directory to store intermediate files, a path to an executable CCFinderX, and information about version control repositories. For each repository, the user is asked to specify the type of version control system (e.g. Git or Mercurial), the root URL of the repository, and the time period that the user is interested in. Table 1 shows example inputs. REPERTOIRE checks the validity of inputs and then proceeds to populate a database with the analysis results of ported edits.

1. Run ‘`src/run_vcs_flow.py`’
  - When an input wizard appears, select “Start a new project”
2. Pick a working directory, e.g. `/var/tmp`.
  - Repertoire creates a sub directory to put its intermediate results.
3. Specify a path to a CCFinder executable.
  - Pick a minimum token size (CCFinder’s input parameter). A minimum token size is the number of lexical token elements that must be similar between two code fragments to be identified as code clones.
4. Select a version control system for each project.
  - REPERTOIRE currently supports Git or Mercurial as target version control systems.
5. Select a URL path for each version control repository.
  - This is the root directory of the repository for Git and Mercurial.
6. Select file extensions for C/C++, headers, and Java files.

7. Select a time period for the project. Repertoire then extracts *diff-based* patches for each commit revision within the time period.
8. Confirm analysis of the given data and then wait for analysis to complete.
9. When the analysis is complete, check the output file created by REPERTOIRE in the working directory.
  - A pickle file called `rep_db.pickle` is generated. Pickle is a file format for Python object serialization and de-serialization. This is used as an input for the visualization and analysis step.

## C. RUNNING REPERTOIRE

1. Run `rep_analysis.py` from `src/analysis`
2. Select the pickle file `rep_db.pickle` produced from the previous step and press **Next**.
3. The GUI provides four analysis views shown in Figure 3: Porting Frequency View, File Distribution View, Developer View, and Porting Latency View (Timing Analysis).

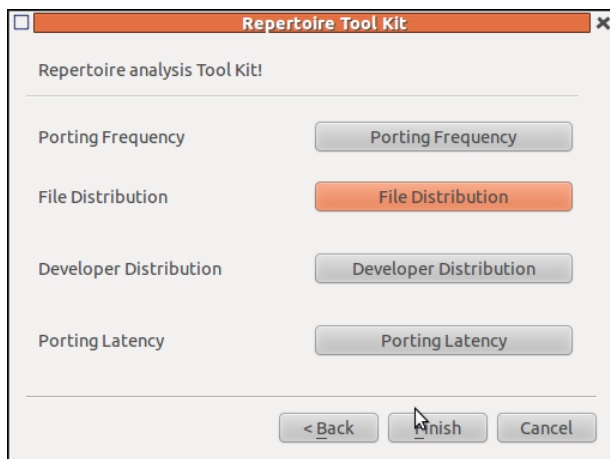


Figure 3: REPERTOIRE Analysis Menu

### C.1 Porting Frequency View

Given the version histories of two projects, this view shows the extent of edits ported from one project to another over the available history. This is represented as a line diagram, where x-axis shows a time line, and the y-axis shows the average percentage of ported edits with respect to total edits in *diff-based* patches. A user may select to see only ported edits from Project A to B, B to A, or both at once. She may also specify a specific time period over the entire evolution history. Steps to run this view:

1. Select **Porting Frequency** in the menu.
2. Select a project: Project A or/and Project B
3. Set a time period for analysis.

### C.2 File Distribution View

This view is a scatter plot where files from Project A is shown along the x-axis and files from Project B is shown along the y-axis. A point is plotted at  $(x,y)$  if there is an edit ported from file X to file Y or vice versa. The color of the dot indicates a ratio of ported edits to total edits. The darker the color is, the higher density of ported edits. Steps to run this view:

1. Select **File Distribution** in the menu.
2. By default, this view does not show full file names. To see the full file names, click **Display Label** option. Alternatively, click on a specific point in the diagram. Corresponding file-names are shown at the bottom.
3. To browse ported code between the selected file pair, press **Display Ported Edit**.
  - A window will show the locations of the ported edits between the two files, along with developer and commit date information.
  - Select any ported code fragment from this list, to browse the corresponding ported edits.

### C.3 Developer View

This view is a scatter plot diagram with developers of project A in x-axis and developers of project B in y-axis. It reflects the interaction pattern of the developers while porting code. Another pie chart shows which developers are responsible for what fraction of the total ported lines in a particular project. Steps to run the developer analysis:

1. Select **Developer Distribution** in the menu.
  - This view shows a scatter plot of developer distribution, i.e., a point is plotted at  $(x,y)$  if developers at x port code written by developer at y, and vice versa.
2. By default, this view does not show developers' identity, as it clutters the display. To see the identities press **Display Label**. Alternatively, click on any point on the scattered plot. Corresponding developers names are shown at the bottom.
3. To see the distribution of the developer who ported code in a particular project in the form of pie chart, select project A and/or project B from right hand window. Then press **Display Developer Porting Statistics**.

### C.4 Porting Latency View

This analysis shows how long it takes for a patch to be propagated to the other project on average. Steps to run this analysis:

1. Select **Porting Latency** in the menu.
2. Select a project (e.g. Project A and/or B), and then press **Porting Latency** button. A cumulative distribution of the time taken to port edits from the source to the target projects is shown when **Cumulative Distribution** is selected.