

Finding Predictors of Field Defects  
for Open Source Software Systems in  
Commonly Available Data Sources:  
a Case Study of OpenBSD

Paul Luo Li, Mary Shaw, Jim Herbsleb,  
June 2005  
CMU-ISRI-05-121

Institute for Software Research International  
School of Computer Science  
Carnegie Mellon University  
Pittsburgh PA, 15213

This paper is an expanded version of the paper titled: Finding Predictors of Field Defects for Open Source Software Systems in Commonly Available Data Sources: a Case Study of OpenBSD, in *METRICS*, 2005.

This research was supported by the National Science Foundation under Grants ITR-0086003, IIS-0414698, and CCF-0438929, by the Carnegie Mellon Sloan Software Center, and by the High Dependability Computing Program from NASA Ames cooperative agreement NCC-2-1298.

Keywords: Process metrics, product metrics, software science, software quality assurance, measurement, documentation, reliability, experimentation, field defect prediction, open source software, reliability modeling, CVS repository, request tracking system, mailing list archives, deployment and usage metrics, software and hardware configurations metrics

## **ABSTRACT**

Open source software systems are important components of many business software applications. Field defect predictions for open source software systems may allow organizations to make informed decisions regarding open source software components. In this paper, we remotely measure and analyze predictors (metrics available before release) mined from established data sources (the code repository and the request tracking system) as well as a novel source of data (mailing list archives) for nine releases of OpenBSD. First, we attempt to predict field defects by extending a software reliability model fitted to development defects. We find this approach to be infeasible, which motivates examining metrics-based field defect prediction. Then, we evaluate 139 predictors using established statistical methods: Kendall's rank correlation, Pearson's rank correlation, and forward AIC model selection. The metrics we collect include product metrics, development metrics, deployment and usage metrics, and software and hardware configurations metrics. We find the number of messages to the technical discussion mailing list during the development period (a deployment and usage metric captured from mailing list archives) to be the best predictor of field defects. Our work identifies predictors of field defects in commonly available data sources for open source software systems and is a step towards metrics-based field defect prediction for quantitatively-based decision making regarding open source software components.

# 1 INTRODUCTION

Open source software systems such as operating systems are important components of many business software applications. Being able to predict field defects (customer reported software problems requiring developer intervention to resolve) may allow existing quantitatively-based decision making methods to be used to:

Help organizations that are seeking to adopt open source software components to make informed choices between candidates

Help organizations using open source software components to decide whether they should adopt the latest release

Help organizations that adopt a release to better manage resources to deal with possible defects

In this paper, we present a case study of the open source operating system OpenBSD in which we try two different approaches to predicting field defects: model fitting and a metrics-based approach.

Prior work by Li et. al. [16] shows that the Weibull model is the preferred model for modeling the defect occurrence pattern of OpenBSD. In the work we report here, we attempt to predict field defects by extending a Weibull model from development to the field. We find that it is not possible to fit an acceptable Weibull model to development defects. The release dates of OpenBSD are consistently around the time when the rate of defect occurrences peaks. Hence, there is insufficient data to fit a Weibull model. This result is consistent with Kenny's finding in [7]. Our finding that it is not possible to fit a Weibull model until the rate of defect occurrences establishes the need for metrics-based field defect prediction.

Identifying and collecting predictors (metrics available before release) are prerequisites activities for metrics-based field defect prediction. We attempt first steps toward a metrics-based field defect prediction model by identifying and collecting potentially important predictors of field defects for OpenBSD. Prior work has identified important predictors of field defects and has predicted field defects for commercial software systems (e.g. Khoshgoftaar et. al. [9], Ostrand et. al. [29], Mockus et. al. [22]). The categories of predictors used in prior work are product metrics, development metrics, deployment and usage (DU) metrics, and software and hardware configurations (SH) metrics. However, prior work has not examined open source software systems, has not examined all categories of predictors simultaneously, and has not identified commonly available data sources for each category of predictor. In this paper, we examine predictors of field defects for an open source software system.

Our experiments show it is possible to collect product, development, DU, and SH predictors from data sources commonly available for open source projects. We identify seven important predictors collected from mailing list archives and the CVS code repository. Somewhat surprisingly, the most important predictor for the OpenBSD project is the number of messages to the technical discussion mailing list during the development period, which is a deployment and usage metric collected from mailing list archives.

Section 2 discusses prior work and motivates our work. Section 3 describes OpenBSD. Section 4 and 5 discuss our data collection method and data analysis method. Section 6 presents the results. Section 7 contains a discussion of our findings. Section 8 is the conclusion.

## 2 PRIOR WORK AND MOTIVATION

In this section, we motivate our work by discussing prior work. We define field defects as user-reported code-related problems requiring programmer intervention to correct. This is the same

definition used by Li et. al. in [16]. We discuss software reliability modeling, software metrics as predictors, and methods used to establish predictors as important.

## **2.1 Software reliability modeling**

Prior work by Li et. al. [16] shows that it is possible to model the rate of field defect occurrences of OpenBSD using the Weibull model. Software reliability modeling research summarized by Lyu in [17] suggests that it may be possible to predict field defects by fitting a Weibull model to development defects and then extending the model to the field. This leads to our first question:

*Is it possible to predict field defects by fitting a Weibull model to development defects and then extending the model to the field?*

Li et. al. use non-linear least squares (NLS) regression to fit defect occurrence data to the Weibull model. We will use the same regression technique to fit a Weibull model to development defects.

## **2.2 Software metrics as predictors**

Metrics are defined by Fenton and Pfleeger in [4] as outputs of measurements, where measurement is defined as the process by which values are assigned to attributes of entities in the real world in such a way as to describe them according to clearly defined rules. Metrics available before release are predictors, which can be used to predict field defects.

We categorize predictors used in prior work using an augmented version of the categorization schemes used by Fenton and Pfleeger in [4], Khoshgoftaar et. al. in [15], and the IEEE standard for software quality metrics methodology [3]:

**Product metrics:** metrics that measure the attributes of any intermediate or final product of the software development process [3]. The product metrics used in prior work are computed using a snapshot of the code. Product metrics have been shown to be important predictors by studies such as Khoshgoftaar et. al. [10], Jones et. al. [6], Khoshgoftaar et. al. [11], Khoshgoftaar et. al. [15], and Khoshgoftaar et. al. [13].

**Development metrics:** metrics that measure attributes of the development process. The development metrics used in prior work are usually computed using information in the change management system or the version control system. Development metrics have been shown to be important predictors by studies such as Khoshgoftaar et. al. [10], Khoshgoftaar et. al. [11], Khoshgoftaar et. al. [15], and Khoshgoftaar et. al. [13].

**Deployment and usage metrics (DU):** metrics that measure attributes of deployment of the software system and usage in the field. Little prior work has examined DU metrics, and no data source is consistently used. DU metrics have been shown to be important predictors by studies such as Jones et. al. [6], Khoshgoftaar et. al. [11], Khoshgoftaar et. al. [15], Khoshgoftaar et. al. [13], and Mockus et. al. [22].

**Software and hardware configurations metrics (SH):** metrics that measure attributes of the software and hardware systems that interact with the software system in the field. Little prior work has examined SH metrics and no data source is consistently used. SH metrics have been shown to be important predictors by Mockus et. al. [22].

Prior work has examined only commercial software systems, and no prior work has examined predictors in all the categories simultaneously. Product and development predictors commonly used in prior work can be computed for open source software systems since the data sources used to compute the predictors (e.g. snapshots of the code and information in the version control system) are commonly available for open source software projects. However, prior work examining DU and SH metrics has had accurate information about system deployment and the users, such as deployment logs and usage profiles (e.g. Khoshgoftaar et. al. [11]) or information from a customer monitoring system (Mockus et. al. [22]). The data sources used in those studies

are generally not available for open source projects; therefore, the DU and SH metrics used in prior work are not available. These concerns lead us to two additional questions:

*Is it possible to collect DU, and SH predictors using data sources commonly available for open source software projects?*

*What are the important predictors of field defects for OpenBSD?*

## **2.3 Methods of establishing the importance of predictors**

Three methods are commonly used to establish a predictor as important:

1. Show high correlation between the predictor and field defects. This method is recommended by IEEE [3] and is used by Ohlsson and Alberg [25] and Ostrand and Weyuker [28].
2. Show that the predictor is selected using a model selection method. This method is used by Jones et al. [6] and Mockus et al. [22].
3. Show that the accuracy of predictions improves with the predictor included in the prediction model. This method is used by Khoshgoftaar et al. [10] and Jones et al. [6].

We use methods 1 and 2 to determine important predictors in this paper. Since we examine predictors that may be included in a metrics-based field defect prediction model but do not actually produce a prediction model, we do not use method 3.

## **3 SYSTEM DESCRIPTION**

In this section, we present the open source software system OpenBSD. We present project details, information on the code repository, information on the request tracking system, and information on the mailing list archives.

### **3.1 Project details**

OpenBSD is a Unix-style operating system written primarily in C. The project dates back to 1995 and has developers (i.e. users who have the write access rights to the CVS code repository) in North America, South America, Europe, Australia, and Asia. This project is similar to the FreeBSD project examined by Dinh-Trong and Bieman [33].

We examine the project between approximately 1998 and 2004. During that time, there were 10 releases (of which we examine 9, as we explain below) and the CVS code repository documented development changes by 159 different developers.

The OpenBSD project uses the Berkley copyrights. The Berkley copyrights retain the rights of the copyright holder, while imposing minimal conditions on the use of the copyrighted material [26]; therefore, OpenBSD has been incorporated into several commercial products.

The OpenBSD project puts out a release approximately every six months. The release dates are published on the web [26].

### **3.2 The code repository**

The OpenBSD project manages its source code using a CVS code repository. Developers are users who have both read and write access rights. Someone becomes a “developer” (i.e. getting an account on the main server) by “doing some good work and showing that he/she can work with the team” [26]. Everyone else has read access to the CVS code repository.

### **3.3 The request tracking system**

The OpenBSD project uses a problem tracking system. Anyone can report a problem by using the *sendbug* command built into OpenBSD [26]. Each problem report is assigned a unique number and stored in the bugs database. The problem report can be tracked on-line using the unique number. A problem report can be assigned one of four classes: sw-bug (software bug), doc-bug

(documentation bug), change request, and support. All problem reports are initially marked as open, and then a developer acts on the report and changes the status accordingly.

Our measure of defects for OpenBSD is user submitted problem reports in the request tracking system of the class software bugs. We count each problem report (which may not be unique) because a user deemed the problem important enough to report. These software related problem reports require a developer's intervention to resolve. This measure of defects is used by Li et. al. in [16]. Defects that occur after the release date are considered field defects. Defects that occur during the development and test period are considered development defects.

### **3.4 Mailing lists**

The OpenBSD project has 23 mailing lists in five categories:

- General interest lists
- Developer's lists
- Platform specific lists
- CVS changes mailing lists
- CTM (emails out deltas to the source).

Not all lists are active and not all lists are archived consistently. The two most complete archives are at sigmasoft [31] and MARC [18].

## **4 DATA COLLECTION**

This section describes the data collection process we used to extract data from the request tracking system, the CVS repository, and the mailing lists. It also describes the predictors we collect.

We consider the published date of release (announced on the OpenBSD website) rounded to the nearest month to be the release date for the release. We round the date to the nearest month due to the time it takes to install the operating system, use the system, and discover and report a problem. Someone reporting a bug right after the un-rounded release date is unlikely to be using the latest release. This is the same approach taken by Mockus et. al. in [22].

We consider the date of the first reported defect rounded to the nearest month to be the start of development. The date of the first reported defect usually occurs several months before the date of release and represents the first time when a problem can be reported against the release that is in development. The development period is then the duration between the start of development and the release date.

### **4.1 The request tracking system**

We wrote Java programs and perl programs to download each problem report from the OpenBSD website and parse the report to extract the report open date, the class (e.g. sw-bug, doc-bug, or change request), the release reported against, and the machine (i.e. the hardware configuration such as i386 or sparc).

There was one anomaly. Three months of data were missing between August 2002 and October 2002. We verified this by examining the bugs mailing list archive (i.e. the mailing lists that records messages to the request tracking system). The mailing list archive showed no bugs recorded during that time interval even though there is activity on the bugs mailing list, which indicates that problems did occur. This happened during development and deployment of release 3.2. As a result, we did not examine release 3.2.

### **4.2 The CVS repository**

We used the CVS checkout command to download the tagged release version of the source code from the CVS repository for releases 2.4 to 3.3 (except release 3.2).

We used four metrics tools and several scripts to compute product metrics from the C source files. The tools we used were:

- RSM by M Squared Technologies [23]
- SourceMonitor by Campwood Software [1]
- c\_count written by Thomas E. Dickey [2]
- metrics written by Brian Renaud [19]

We arrived at these tools by conducting a web search, asking experts for help, and posting to the comp.software-eng and comp.software.measurement newsgroups. We evaluated the collection of tools and selected those listed above.

We encountered an existing CVS bug when downloading the source code for release 2.4 and release 2.5. As a result, we had to bypass a directory that contained HTML help documents. We also encountered 10 files with coding anomalies that the metrics tools could not resolve. We skipped those files for all releases. These files constitute less than .1% of the number of C source files.

We used the CVS log command to obtain information on changes to the source code. We used the log information between the start of development of release 2.4 and the release date of release 3.3. There were 97,566 committed changes in the development periods of the nine releases.

### **4.3 Mailing list archives**

We wrote java programs to extract the number of messages posted each month in the mailing lists archives.

Not all lists were archived consistently and not all lists were active. Consistent data was not available for many of the lists before 1998; therefore, we only considered releases after 1998. When an archive showed no messages for a certain month, we were often unable to determine if no messages were posted or if the archive failed to properly record messages (both of which occur). Therefore, lists that had intervals in which no messages were posted for more than three months were not considered.

### **4.4 Metrics**

We provide a summary of the 139 predictors we collected. Appendix B lists the full set of metrics and Appendix C contains the collected data.

#### **4.4.1 Product metrics**

We collected 101 product metrics using snapshots of the code from the CVS code repository (details in Appendix B and C). Due to tooling constraints, we did not collect all the product metrics used in the literature. However, we did collect metrics that covered all of the dimensions of variation in the product metrics identified by prior work. Munson and Khoshgoftaar identified the dimensions of product metrics (i.e. components of variance captured by product predictors) used in the literature using principal component analysis in [24]. Principal component analysis captures the dimensions of variance in a group of predictors. Predictors that load on the same principal component capture the same dimension of variance and are highly correlated with each other [24]; therefore, it may be sufficient to use a predictor from each dimension. We give the dimensions and examples of the product metrics we used to capture the variation in the dimension in table 1. The product metrics we used had been shown to load on the principal component by Munson and Khoshgoftaar in [24]. Appendix B lists the full set of product metrics.



**Table 1. Product metrics**

<i>Dimension</i>	<i>Product metrics used in this study</i>
Control: metrics related the flow of program control	<i>Cyclo</i> : Cyclomatic complexity <i>KWbreak</i> : Number of occurrences of the key word break (which is equivalent to possible program knot count as shown by Khoshgoftaar and Szabo [14])
Action: number of distinct operations and statements	<i>UOpand</i> : Unique operands <i>UOpator</i> : Unique operators
Size: size or item count of a program	<i>Statements</i> : Total number of statements per file summed across all files <i>LOC</i> : Lines of code per file summed across all files
Effort: Halstead’s effort metrics	<i>PGeffort</i> : Halstead’s effort metric per file summed across all files
Modularity: degree of modularity of a program	<i>DeepNest</i> : Number of statements at nesting level >9 per file summed across all files

#### 4.4.2 Development metrics

We collected 22 development metrics (details in Appendix B and C). Due to differences in the style of development, we were not able to collect the same development metrics used in the literature. However, we tried to collect metrics that captured the same intent as the metrics used in the literature in our study. We collected metrics that cover all of the independent dimensions of variation in the development metrics identified by Khoshgoftaar et. al. in [13] and [15]. Khoshgoftaar et. al. used principal component analysis to identify the dimensions of variation in their development metrics in [13] and [15]. Khoshgoftaar et. al. examined a commercial software system while we examined an open source software system; therefore, we made changes to the metrics to account for the differences between commercial and open source software systems. We offer an interpretation of the dimensions captured by each principal component (which is not offered by Khoshgoftaar et. al.), examples of the metrics belonging to each dimension in [13] and [15], and the metrics we used to capture the same sources of variance in table 2. We made one major modification. Since OpenBSD did not distinguish between designers and testers, we combined the dimensions identified by Khoshgoftaar et. al. that separated designers and testers. We believe our metrics captured the same source of variation as the referenced metrics since the only changes we made were to accommodate differences between commercial and open source styles of development. (Metrics collected using the CVS code repository are indicated by ‘CVS’, ones collected using the request tracking system are indicated by ‘RTS’, and ones collected using mailing list archives are indicated by ‘MLA’.)

**Table 2. Development metrics**

<i>Dimensions</i> [13] and [15]	<i>Example of metrics in dimensions</i> [13]	<i>Development metrics used in this study</i>
Dimension 1: the number of changes	Total number of changes to the code for any reason	<i>TotalUpdate (CVS)</i> : Total number of updates during the development period

Dimension 2: experience of the people making changes	Number of updates to this module by designers who had 10 or less total updates in entire company career	<i>BotHalfC</i> (CVS): Number of different developers making changes to files that are c source files during the development period who are in the bottom 50% of all developers ranked based on the number of changes
Dimension 3: amount of change to the code	Net increase in lines of code	<i>Difference</i> (CVS): Lines added to c source files minus lines deleted from c source files during the development period
Dimension 4 and 7: problems found during the development of the prior release	Number of problems fixed that were found by designers or during beta testing in the prior release	<i>PreBugsPrev</i> (RTS): Total number of field defects reported during the development period of the previous release
Dimension 5 : field problems found by customers in prior releases	Number of problems fixed that were found by customer in the prior release	<i>PreBugsAll</i> (RTS): Total number of field defects reported during the development period in all releases
Dimension 6 and 8: problems found during the development of the current release	Number of problems found by designers or during beta testing in the current release	<i>PreBugsCurrent</i> (RTS): Number of field defects reported against the release under development during the development period

#### 4.4.3 Deployment and usage metrics

We collected nine deployment and usage metrics (details in Appendix B and C). The metrics we collected fall into two categories: mailing list predictors and request tracking system predictors. Mailing list predictors counted the number of messages to non-hardware related mailing lists during development. We believed our mailing list predictors captured characteristics of deployment and usage because they quantified the amount of interest in OpenBSD, which might be related to how many systems were deployed and how much the systems were used. Request tracking predictors counted the number of problem reports during development that were not defects. We believed our request tracking system predictors captured characteristics of deployment and usage because users had to install OpenBSD and use the system before they could report a problem. We present the two categories, examples of predictors in the categories, and short justifications for the predictors in table 3. Appendix B lists the full set of deployment and usage metrics.

Table 3. Deployment and usage metrics

<i>Category of predictors</i>	<i>DU metrics used in this study</i>	<i>Justification</i>
Mailing list predictors	<i>MiscMailings</i> (MLA): number of messages to the miscellaneous mailing list, a general interest mailing list, during the development period <i>AdvocayMailings</i> (MLA): number of messages to the advocacy mailing list (which promotes the use of OpenBSD), a general interest mailing list,	These metrics quantify the amount of interest in OpenBSD, which may be related to how many systems are deployed and how much the systems are used.

	during the development period	
Request tracking system predictors	<i>ChangeRequests</i> (RTS): Number of change requests during the development period <i>DocBugs</i> (RTS): Number of reported documentation problems during the development period	These metrics quantify the amount of deployment and usage because users must install OpenBSD and use the system before they can request changes or report documentation problems

#### 4.4.4 Software and hardware configurations metrics

We collected seven software and hardware configurations metrics in all (details in Appendix B and C). The metrics we collected fall into two categories: mailing list predictors and request tracking system predictors. Mailing list predictors counted the number of messages to hardware specific mailing lists during development. We believed our mailing list predictors captured characteristics of software and hardware configurations because they reflect the amount of interest/activity related to the specific hardware, which might be related to how many of the specified hardware machines had OpenBSD installed. Request tracking predictors counted the number of defects (field defects and development defects) during development that identify the type of hardware used. We believed our request tracking system predictors captured characteristics of software and hardware configurations because users had to install OpenBSD on the specified HW before they could report a problem. We present the two categories, examples of predictors in the categories, and short justifications of the predictors in table 4. Appendix B lists the full set of software and hardware configurations metrics.

Table 4. Software and hardware configurations metrics

<i>Category of predictors</i>	<i>SH metrics used in this study</i>	<i>Justification</i>
Mailing list predictors	<i>SparcMailing</i> (MLA) Number of messages to the sparc hardware specific mailing list, a platform specific mailing list, during the development period	This metrics may reflect the amount of interest/activity related to the specific hardware, which may be related to how many of the specific hardware machines have OpenBSD installed.
Request tracking system predictors	<i>CurrentBSD Bugs i386HW</i> (RTS): Number of field defects reported against the current release during the development period that identify the machine as type i386	These metrics may quantify the number of machines with specific HW that have OpenBSD installed since users must install and use the system to report a problem

## 5 DATA ANALYSIS

First, we attempted to fit Weibull models to development defects. We used NLS regression to fit the Weibull models. NLS is a widely used model fitting method discussed in detail by Lyu in [17].

Next, we computed the correlations between the predictors and field defects in order to identify important predictors. We did not consider predictors that did not vary since they cannot predict field defects (e.g. we discarded the predictor measuring the number of instances of the key work 'struct' in the code, which was zero for all releases). We computed Spearman's rank correlation

(?), Kendall's rank correlation ( $\tau$ ), and the statistical significance of the correlations. These are standard ways of computing rank correlation. Holland and Wolfe [5] recommended using rank correlation when the data are not normally distributed. We determined that the data were not normally distributed by examining data plots. Refer to Weisberg [35], Venable and Ripley [34], and Holland and Wolfe [5] for detailed explanations of rank correlation.

Finally, we performed a forward AIC model selection to identify important predictors. The predictors selected using the forward AIC model selection method complement each other since each predictor improves the fit substantially (i.e. enough to overcome the AIC penalty) even with the other predictors already in the model. The forward AIC model selection method can be used to select a subset of predictors as a first step in a regression analysis even if the data is not normally distributed. Refer to Weisberg [35] for a detailed explanation. The model selection process usually continues until the AIC score does not improve with additional predictors; however, since we had 9 observations and 139 predictors, we stopped at three iterations to prevent over fitting. Similar model selection methods were used by Ostrand et. al. in [29] and Khoshgoftaar et. al. in [12].

For all our analysis, we used the open source statistical program R [30].

## 6 RESULTS

We present the results of fitting the Weibull model, evaluating the predictors using correlation, evaluating the predictors using forward AIC model selection, and comparing important predictors. We find that the number of messages to the technical mailing list during development is the best predictor.

### 6.1 Prediction using a software reliability model

We are not able to fit a Weibull model to development defects. The NLS model fitting procedure does not converge for any of the releases. Our finding that the modeling fitting procedure does not converge is consistent with Kenny's findings in [7], which show that it is not possible to fit a Weibull model until most of the defects have occurred (i.e. past the hump in the number of defects). A typical release with the release date indicated is shown in figure 1. Plots of all the releases are in Appendix A. In 7 out of 9 releases, the release date is within two months of the time the rate of defect occurrences peak. In 8 out of 9 releases, the release date is either within one month or before the time the rate of defect occurrences peak. We cannot predict field defects by fitting a Weibull model to development defects.

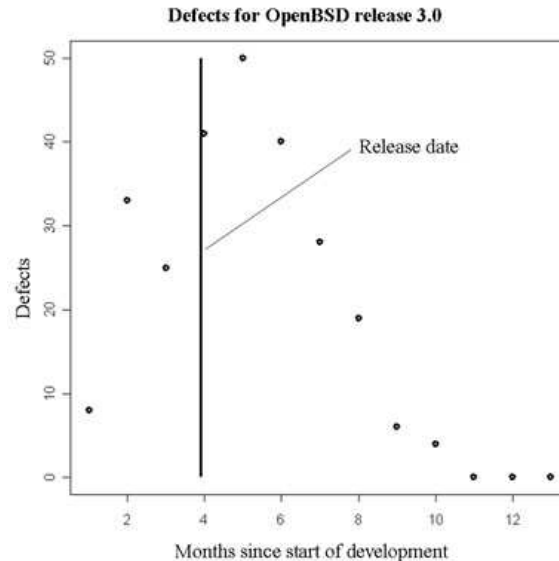


Figure 1. Defects for OpenBSD release 3.0

## 6.2 Analysis of predictors using correlations

Table 5 presents predictors that are significant at the 95% confidence level (CL) using rank correlation (a blank indicates that a predictor’s correlation is not significant at the 95% CL). We briefly explain the predictors in this section.

### Product metrics (computed using a snapshot of the code from the CVS code repository and the RSM metrics tool):

- *TotMeth*: Total number of methods.
- *PubMeth*: Number of public methods.
- *InlineComment*: Number of inline comments.
- *ProtMeth*: Number of protected methods.
- *CommentsClass*: Number of comments in classes summed across all classes.
- *InterfaceCompClass*: Number of parameters + number of returns in classes summed across all classes.
- *TotalParamClass*: Total number of parameters in classes summed across all classes.

### Development metric (computed using the CVS code repository):

- *UpdateNotCFiles*: During the development period, the number of updates (deltas) to files that are not c source files.

### Deployment and usage metric (computed using mailing list archives):

- *TechMailing*: Number of messages to the technical mailing list, a developer’s mailing list, during development.

### Software and hardware configuration metric (computed using mailing list archives):

- *SparcMailing*: number of messages to the sparc hardware specific mailing list, a platform specific mailing list, during the development period.

Table 5. Rank correlations

<i>Predictor</i>	<i>Kendall Correlation</i>	<i>p-value</i>	<i>Spearman Correlation</i>	<i>p-value</i>
TechMailing	0.61	0.02	0.78	0.02
TotMeth	0.61	0.02	0.73	0.03
PubMeth	0.61	0.02	0.73	0.03
CommentsClass	0.61	0.02	-	-
ProtMeth	0.57	0.03	0.67	0.05
InlineComment	0.56	0.04	0.68	0.05
InterfaceCompClass	0.51	0.05	-	-
TotalParamClass	0.51	0.05	-	-

### 6.3 Analysis of predictors using forward AIC model selection

We use three iterations of the forward AIC model selection method to select important predictors. Due to space limitation, we present the final linear model in table 6. The predictors are listed in the order selected. The AIC score of the final model is 75.52. The  $r^2$  between the fitted model and field defects is 0.93. This high correlation suggests possible over fitting and confirms the need to stop at three iterations.

Table 6. AIC selected model

<i>Variable</i>	<i>Estimate coefficient</i>	<i>Standard Error</i>	<i>t value</i>	<i>Pr(&gt; t )</i>
(Intercept)	134.32	18.06	7.437	0.0007
TechMailing	0.1102	0.015	7.445	0.0007
UpdatesNotCFiles	-0.0289	0.005	-5.757	0.0022
SparcMailing	0.1406	0.045	3.153	0.0253

The linear model in table 6 is not intended to be a valid prediction model. Additional steps need to be taken (e.g. adjust for non-constant variance) before the model can be used for prediction. Further validation of the predictors is also needed. We hope to do so in future work.

The estimated coefficients require interpretation. Since ranges of the predictors differ and all predictors are statistically significant, it is sensible to examine only the direction of the estimated coefficients (i.e. if they positive or negative). The coefficient for TechMailing is positive, indicating that increases in the metric correspond to more field defects. TechMailing measures the amount of deployment and usage of the system. This metric quantifies the amount of interest in OpenBSD, which may be related to how many systems are deployed and how much the systems are used. Our finding that increased deployment and usage correspond to more field defects is consistent with findings in Jones et al. [6] and Mockus et al. [22].

The coefficient for UpdatesNotCFiles is negative, indicating that increases correspond to fewer field defects. We think larger UpdatesNotCFiles may indicate maintenance (i.e. efforts to eliminate problems); therefore, it corresponds to fewer field defects. The coefficient for SparcMailing is positive indicating that increases in SparcMailing correspond to more field problems. Increase in SparcMailing may indicate increased activity/usage related to the sparc hardware, which may lead to field defects unaccounted for by the other predictors.

### 6.4 Comparison of important predictors

We compare the important predictors by examining the rank correlation among the predictors and field defects. This may allow us to determine which predictors may produce better predictions. We do not have enough observations to perform a principal component analysis.

The correlations between important predictors selected using rank correlation and field defects in table 7 indicate that increases in each of the predictors correspond to more field defects. These correlations are consistent with findings in prior work. The relationship between TechMailing (a DU metrics) and field defects is consistent with findings in Jones et. al. [6] and Mockus et. al. [22]. All other important predictors are product metrics. Our finding that increases in the product correspond to more field defects is consistent with findings in Ostand et. al.[29] and Jones et. al. [6]. However, the predictors are highly correlated with each other. This suggests that it may be sufficient to use just one of the predictors and that including all the predictors in a model may result in the multi-co-linearity problem discussed in Fetton and Pflieger [4] .

**Table 7. Correlations among important predictors**

		AIC selected predictors			Correlation selected predictors			
	<i>Field defects</i>	<i>Sparc Mailing</i>	<i>Updates NotCFiles</i>	<i>Tech Mailing</i>	<i>Tot Meth</i>	<i>Pub Meth</i>	<i>Inline Comment</i>	<i>Prot Meth</i>
<i>Field defects</i>	1.000	0.278	-0.111	0.611	0.611	0.611	0.556	0.567
<i>SparcMailing</i>	0.278	1.000	0.500	0.111	0.556	0.556	0.278	0.433
<i>UpdatesNotCFiles</i>	-0.111	0.500	1.000	0.167	0.278	0.278	0.222	0.367
<i>TechMailing</i>	0.611	0.111	0.167	1.000	0.444	0.444	0.611	0.500
<i>TotMeth</i>	0.611	0.556	0.278	0.444	1.000	1.000	0.722	0.767
<i>PubMeth</i>	0.611	0.556	0.278	0.444	1.000	1.000	0.722	0.767
<i>InlineComment</i>	0.556	0.278	0.222	0.611	0.722	0.722	1.000	0.833
<i>ProtMeth</i>	0.567	0.433	0.367	0.500	0.767	0.767	0.833	1.000

The correlation among important predictors selected using the forward AIC model selection method are lower than the correlation among important predictors selected using correlations. This confirms that the each predictor selected using model selection captures information not captured by the other predictors; therefore, they will complement each other in a prediction model and avoid the multi-co-linearity problem.

## 7 Discussion

We have established that it is not possible to fit a Weibull model to development defects for OpenBSD. We present results from fitting the Weibull model because prior work has identified the Weibull model as the preferred model. In addition, we also have results from experiments showing that it is not possible to make meaningful field defect predictions by extending other software reliability models fitted to development defects (i.e. the Gamma model, the Logarithmic model, the Exponential model, the Power model). Due to space limitations, those results are omitted. These results motivate the need to consider metrics-based field defect prediction.

We find that it is possible to collect product, development, DU, and SH predictors using data sources commonly available for open source software systems. In addition to validating the CVS code repository and the request tracking system as sources of predictors, we establish mailing list archives as an important data source, one not considered by previous studies.

We find that the most important predictor for the OpenBSD project is TechMailing collected from mailing list archives. The TechMailing predictor is the most highly rank correlated predictor with the number of defects and is the first variable selected using AIC forward model selection. We have validated this finding by talking with developers on the discussion forum. Feedback [27] indicates that this finding fits with the developers’ intuition that participation by active developers (reflected by postings to the TechMailing list) leads to more defect discoveries. A plot of TechMailing against field defects is shown in figure 3. Other important predictors selected include four product metrics collected from the CVS code repository, a development metric collected from the CVS code repository (UpdatesNotCFiles), and a software and hardware configurations metric collected from mailing list archives (SparcMailings).

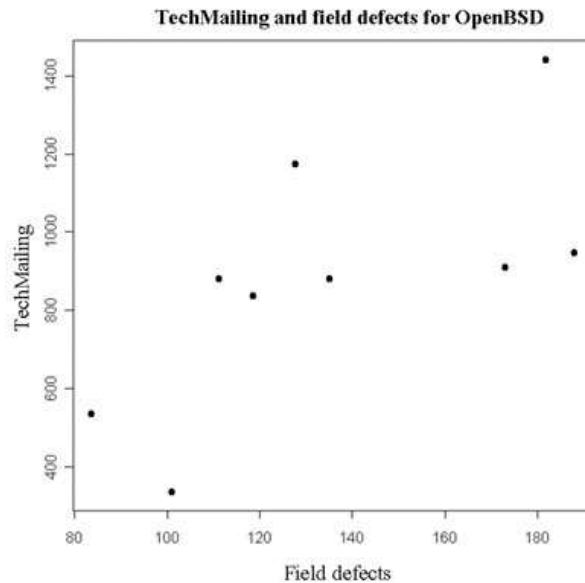


Figure 3. TechMailing and field defects

In contrast to findings in commercial software systems, (e.g. Mockus et. al. [21], Khoshgoftaar et. al. [10], and Khoshgoftaar et. al. [11]) predictors regarding changes to source files and those regarding developers are not important predictors for OpenBSD. We suspect this is due to the review and check-in process employed by the OpenBSD project (and possibly by other open source projects as well), which assures that all changes are of a certain quality regardless of the number of changes or the author of the change. All changes must be checked-in by a developer, who is someone that has shown ability to work on the code. This is supported by the explanation on the project webpage, which details how someone becomes a developer and gains the ability to check-in code. In addition, many changes are reviewed. We find evidence of this by observing logs of committed changes. Many logs contain markers (of the type “developer id” followed by the @ sign, e.g. art@) indicating that another developer has reviewed the changes.

## 8 CONCLUSION

In our case study of OpenBSD, we find that it is not possible to predict field defects by extending a Weibull model fitted to development defects. This indicates the importance of metrics-based field defect prediction models for open source software systems. We also find that it is possible to collect product, development, DU, and SH metrics using commonly available data sources for opens source projects. In addition, we identify important predictors that can be used to construct a field defect prediction model for OpenBSD using modeling methods in the literature. Such a



model can help organizations make informed decisions regarding open source software components.

The paper presents novel and interesting findings, which are appropriate for a case study. However, our experiments need to be replicated on other open source projects. Replications can help verify that the relationships we have established are not due to chance alone. Future studies can include similar projects developing operating systems like FreeBSD or Debian and other types of systems like MySQL or JakartaTomcat.

Replication of our experiments is relatively straightforward since data sources we use are commonly available for open source software systems. For example, all projects hosted by SourceForge [32] use a CVS code repository, a request tracking system, and have mailing lists.

We do not consider non-c source files in our analysis (e.g. perl files and assembly files). These files may contain valuable information. However, since most of the system is written in c, we feel c source files are the most appropriate files to analyze. In release 3.4 (the most recent release we examine), there are ~36384 files in total. Approximately 17578 are c source files, 2378 are perl source files, and 1624 are assembly files. The remaining files are mostly documentation, configuration, and installation files.

There maybe other metrics we have failed to collect. For example, it may be possible to parse the defect reports for more detailed information regarding bugs, such as which software applications were running when the bug occurred. Since the data sources are available to everyone, we encourage others to explore other predictors.

Results in this paper represent a promising step towards quantitatively-based decision making regarding open source software components. The next step is to use the results in this paper and metrics-based modeling methods in the literature to construct metrics-based field defect prediction models and then to compare their predictions (e.g. the trees based method used by Khoshgotaar et. al. in [13], the neural networks method used by Khoshgoftaar et. al. in [14], and the linear regression used by Mockus et. al. in [22]).

## 9 ACKNOWLEDGMENTS

This research was supported by the National Science Foundation under Grants ITR-0086003, IIS-0414698, and CCF-0438929, by the Carnegie Mellon Sloan Software Center, and by the High Dependability Computing Program from NASA Ames cooperative agreement NCC-2-1298. We thank the developers of OpenBSD for answering our postings. We thank the tool vendor who gave us trial licenses. We thank George Fairbanks, Patrick Riley, and Greg Wolfson for their help and insight.

## 10 REFERENCES

- [1] Campwood Software. <http://www.campwoodsw.com/>
- [2] c\_count. [http://dickey.his.com/c\\_count/c\\_count.html](http://dickey.his.com/c_count/c_count.html)
- [3] IEEE standard for a software quality metrics methodology. In *IEEE Std 1061-1998*, 1998.
- [4] Norman Fenton and Shari Pfleeger. *Software Metrics - A Rigorous and Practical Approach*. Chapman & Hall, London, 1997
- [5] Myles Hollander and Douglas A. Wolfe. *Nonparametric statistical inference*. Wiley & Sons, 1973.
- [6] Wendell Jones, John P. Hudepohl, Taghi M. Khoshgoftaar, and Edward B. Allen. Application of a Usage Profile in Software Quality Models. In *3rd European Conference on Software Maintenance and Reengineering*, 1999.
- [7] Garrison Kenny. Estimating defects in commercial software during operational use. In *IEEE TR on Reliability*, 1993.

- [8] Taghi M. Khoshgoftaar, Edward B. Allen, Wendell Jones, and John Hudepohl. Which software modules will have faults that will be discovered by customers? In *Journal of Software Maintenance: Research and Practice*, 1999
- [9] Taghi M. Khoshgoftaar, Edward B. Allen, Kalai S. Kalaichelvan, and Nitith Goel. Predictive modeling of software quality for very large telecommunications systems. In *IEEE International Conference on Communications*, 1996.
- [10] Taghi M. Khoshgoftaar, Edward B. Allen, Kalai S. Kalaichelvan, Nitith Goel, John Hedepohl, and Jean Mayrand. Detection of fault-prone program modules in a very large telecommunications system. In *Proceedings of ISSRE*, 1995.
- [11] Taghi M. Khoshgoftaar, Edward B. Allen, Xiaojing Yuan, Wendell D. Jones, and John P. Hudepohl. Preparing measurements of legacy software for predicting operational faults. In *ICSM*, 1999.
- [12] Taghi Khoshgoftaar, Adhijit Pandya, and David Lanning. Application of neural networks for predicting program faults. In *Annals of Software Engineering*, 1995.
- [13] Taghi M. Khoshgoftaar, Ruqun Shan, and Edward B. Allen. Using product, process, and execution metrics to predict fault-prone software modules with classification trees. In *HASE*, 2000.
- [14] Taghi Khoshgoftaar and Robert Szabo. Using neural networks to predict software faults during testing. In *IEEE Transaction on Reliability*, 1996.
- [15] Taghi M. Khoshgoftaar and Vishal Thaker and Edward Allen. Modeling fault-prone modules of subsystems. In *Proceedings of ISSRE*, 2000.
- [16] Paul Luo Li, Mary Shaw, Jim Herbsleb, Bonnie Ray, and P. Santhanam. Empirical Evaluation of Defect Projection Models for Widely-deployed Production Software Systems. *FSE*, 2004.
- [17] Michael Lyu. *Handbook of Software Reliability Engineering*. IEEE Society Press, USA, 1996.
- [18] MARC. <http://marc.theaimsgroup.com/>
- [19] metrics. <http://www.chris-lott.org/resources/cmetrics/>
- [20] Audris Mockus, Roy Fielding, and James Herbsleb. A case study of open source software development: the Apache server. *ICSE*, 2000.
- [21] Audris Mockus, David Weiss, and Ping Zhang. Understanding and predicting effort in software projects. In *ICSE*, 2003
- [22] Audris Mockus and Ping Zhang and Paul Luo Li. Drivers for Customer Perceived Quality. In *ICSE*, 2005.
- [23] M Squared Technologies. <http://msquaredtechnologies.com>
- [24] John Munson and Taghi Khoshgoftaar. The dimensionality of program complexity. In *ICSE*, 1989.
- [25] Niclas Ohlsson and Hans Alberg. Predicting fault-prone software modules in telephone switches. In *IEEE Transactions on Software Engineering*, 1996
- [26] OpenBSD. [www.openbsd.org](http://www.openbsd.org).
- [27] OpenBSD discussion thread. <http://marc.theaimsgroup.com/?t=110788031900001&r=1&w=2>
- [28] Thomas J. Ostrand and Elaine J. Weyuker. The Distribution of Faults in a Large Industrial Software System. In *ISSTA*, 2002.
- [29] Thomas J. Ostrand and Elaine J. Weyuker and Thomas Robert M. Bell. Where the bugs are. In *ISSTA*, 2004.
- [30] R. <http://www.r-project.org/>
- [31] Sigmasoft. <http://www.sigmasoft.com/~openbsd/>
- [32] SourceForge. <http://sourceforge.net/>
- [33] Trimg Dinh-Trong and James M. Bieman, Open source software development: a case study of FreeBSD. *Metrics*, 2004.
- [34] W.N. Venables and Brian D. Ripley. *Modern Applied Statistics with S-plus, 4<sup>th</sup> edition*. Springer-Verlag, 2000.
- [35] Sanford Weisberg. *Applied Linear Regression, 2<sup>nd</sup> Edition*. Wiley and Son, 1985.
- [36] Xiaohong Yuan, Taghi Khoshgoftaar, Edward Allen, and K Gasesan. An application of fuzzy clustering to software quality prediction. In *IEEE Symposium on Application-Specific Systems and Software Engineering Technology*, 2000.

# 11 APPENDIX A

In this section, we present results of fitting software reliability models to development defects. We consider the same set of models considered in Li et al. [16].

The Weibull model, Exponential model, and Logarithmic model could not be fitted using data between first availability and general availability for any of the releases. The Gamma model could be fitted for only 2 of the releases (R2.8 and R3.3). However, the predictions are not accurate. The predicted total number of defect occurrences are 71.5 ~ 72 occurrences and 169.9 ~ 170 occurrences for the two releases respectively. The actual total numbers of defect occurrences for the releases are 239 and 172. The power model could be fitted for all the releases, but the fitted models were not meaningful. Each model was strictly increasing. The fitted models and the actual defect occurrences for OpenBSD 2.8 and 3.3 are in figures A1-A2.

We conclude that it is not possible to fit a model using data between first availability and general availability to predict the number of defect occurrences after general availability. This agrees with Kenny's findings in [7] that it is not possible to fit Weibull model parameters to a model until most of the defects have occurred. The idea is that for initially increasing and eventually decreasing defect occurrence patterns, most of the defects need to have already occurred (i.e. past the hump) before a software reliability model can be fitted. All of the releases with the release dates labeled are in figures A3-A11.

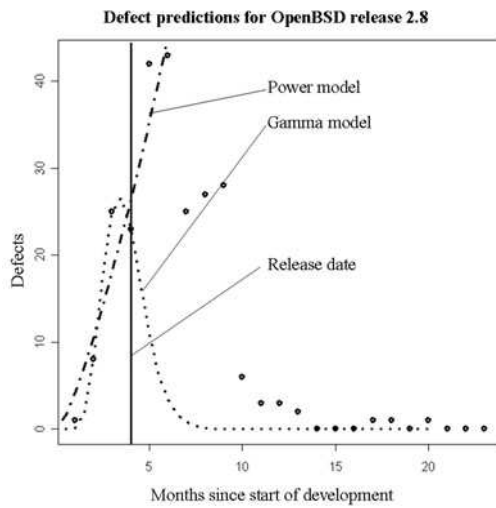


Figure A1. Field defects and fits release 2.8

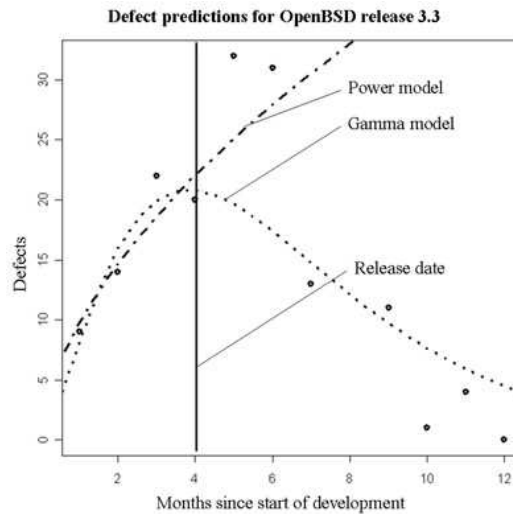
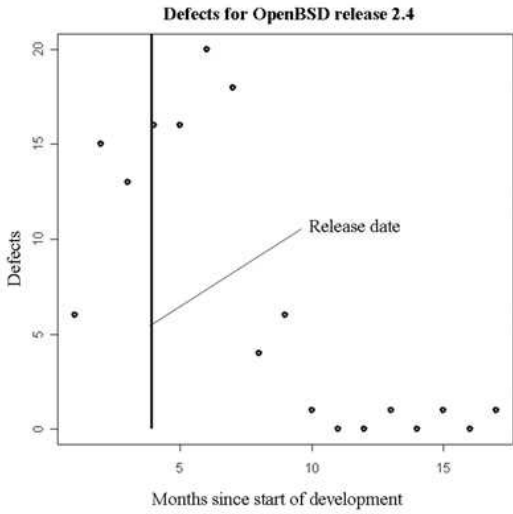
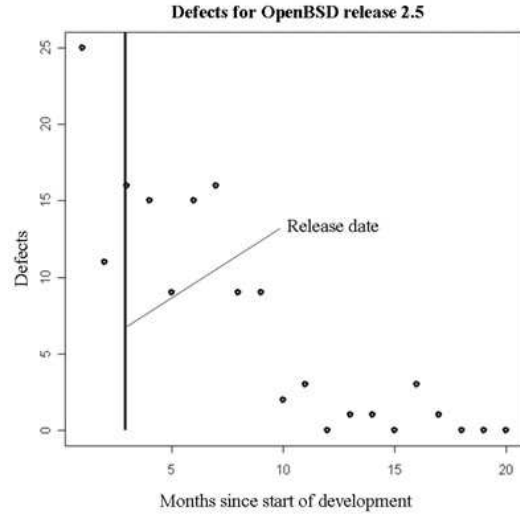


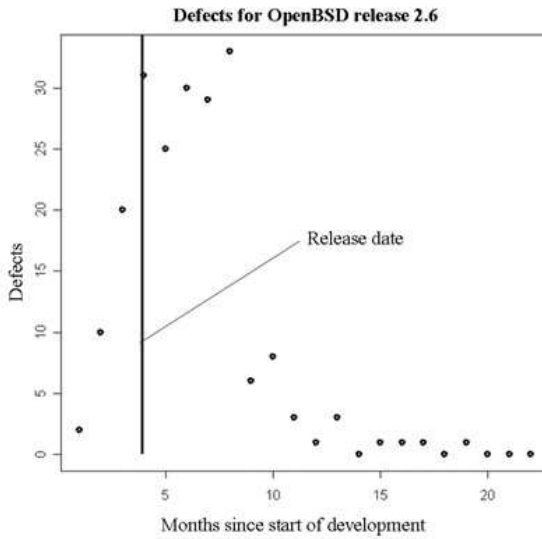
Figure A2. Field defects and fits release 3.3



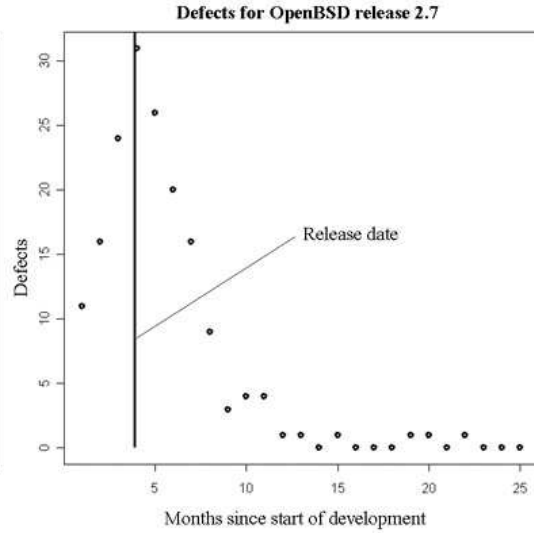
**Figure A3. Field defects release 2.4**



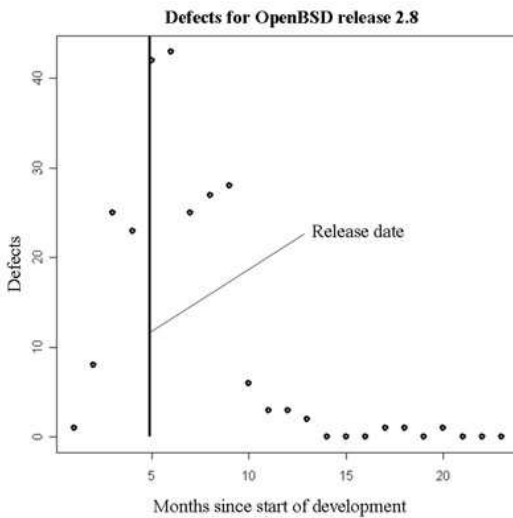
**Figure A4. Field defects release 2.5**



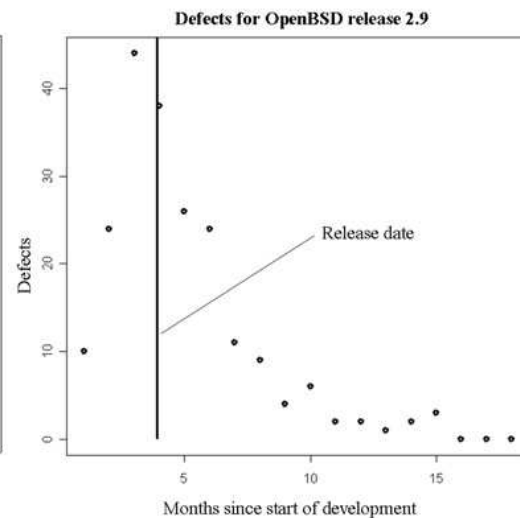
**Figure A5. Field defects release 2.6**



**Figure A6. Field defects release 2.7**



**Figure A7. Field defects release 2.8**



**Figure A8. Field defects release 2.9**

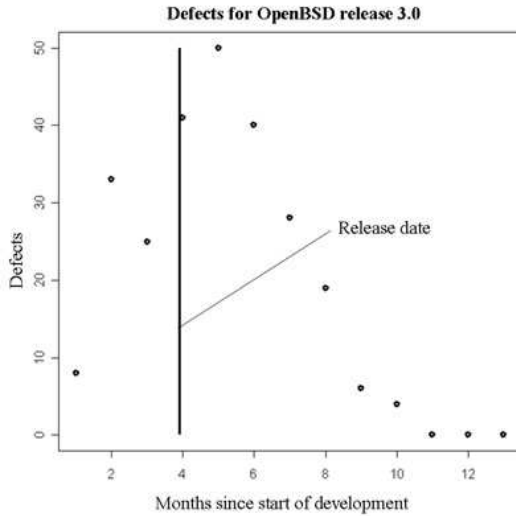


Figure A9. Field defects release 3.0

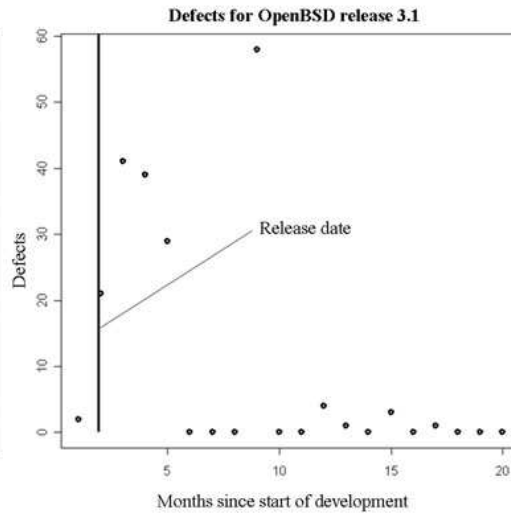


Figure A10. Field defects release 3.1

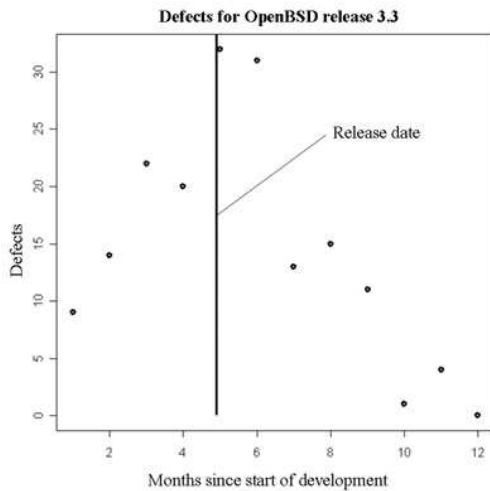


Figure A11. Field defects release 3.3

## 12 APPENDIX B

In this section, we present the full set of metrics used in our study.

### 12.1 Product metrics

*LOC*: Lines of code calculated by the metrics tool Source Monitor

*Statements*: Statements in C and .h files calculated by the metrics tool Source Monitor

*Functions*: Statements calculated by the metrics tool Source Monitor

*Bandwidth*: Modified Bandwidth metric calculated using statements and nesting depth information from the metrics tool Source Monitor. Source Monitor only count nesting up to 10 levels. Therefore, the metrics clip the statements at nesting of 10 levels.

*DeepNest*: Statements at nesting level greater than 10 calculated using the metrics tool Source Monitor

*Cstatements*: Statements in C source files calculated by the metrics tool C\_Count

*Lineswithcomments*: Lines with comments in C source files calculated by the metrics tool C\_Count

*InlineComment*: Lines with inline comments in C source files calculated by the metrics tool C\_Count

*Blanklines*: Blank lines in C source files calculated by the metrics tool C\_Count

*Linesforpreprocess*: Pre-processor lines in C source files calculated by the metrics tool C\_Count

*Lineswithcode*: Lines with code in C source files calculated by the metrics tool C\_Count

*Totallines*: Total number of lines in C source files calculated by the metrics tool C\_Count

*Statements*: Total number of statements in C source files calculated by the metrics tool C\_Count

*Commentchars*: Total number of comment characters in C source files calculated by the metrics tool C\_Count

*Nontextcommentchars*: Total number of non-text comment characters in C source files calculated by the metrics tool C\_Count

*Whitespacechars*: White space characters in C source files calculated by the metrics tool C\_Count

*Preprocessorchars*: Characters for the pre-processor in C source files calculated by the metrics tool C\_Count

*Statementchars*: Statement characters in C source files calculated by the metrics tool C\_Count

*Totalchars*: Total number of characters in C source files calculated by the metrics tool C\_Count

*Tokens*: Total number of recognized program token in C source files calculated by the metrics tool Metrics

*UOpand*: Unique operands in C source files calculated by the metrics tool Metrics

*UOpator*: Unique operators in C source files calculated by the metrics tool Metrics

*TPpand*: Total operands in C source files calculated by the metrics tool Metrics

*TOpator*: Total operators in C source files calculated by the metrics tool Metrics

*PGlength*: Halstead's estimated program length in C source files calculated by the metrics tool Metrics

*PGvolume*: Halstead's program volume in C source files calculated by the metrics tool Metrics

*PGeffort*: Halstead's estimated program effort in C source files calculated by the metrics tool Metrics

*PGlevel*: Halstead's estimated program level in C source files calculated by the metrics tool Metrics

*Files*: Number of C files calculated by the metrics tool RSM

*FunctionsRSM*: Number of functions calculated by the metrics tool RSM

*PhysicalLines*: Number of physical lines calculated by the metrics tool RSM

*LOC*: Lines of code calculated by the metrics tool RSM

*eLOC*: Effective lines of code calculated by the metrics tool RSM

*lLOC*: Logical lines of code calculated by the metrics tool RSM

*Cyclo*: Cyclomatic complexity calculated by the metrics tool RSM

*InterfaceComp*: Interface complexity calculated by the metrics tool RSM

*TotalParams*: Total parameters used calculated by the metrics tool RSM

*TotalReturn*: Total returns calculated by the metrics tool RSM  
*CommentsRSM*: Number of comments calculated by the metrics tool RSM  
*BlanksRSM*: Number of blank lines calculated by the metrics tool RSM  
*Classes*: Number of classes calculated by the metrics tool RSM  
*NestedClasses*: Number of nested classes calculated by the metrics tool RSM  
*TotMeth*: Number of methods in classes calculated by the metrics tool RSM  
*PubMeth*: Number of public methods in classes calculated by the metrics tool RSM  
*Publicattributes*: Number of public attributes in classes calculated by the metrics tool RSM  
*ProtMeth*: Number of protected methods in classes calculated by the metrics tool RSM  
*Protectedattributes*: Number of protected attributes in classes calculated by the metrics tool RSM  
*PrivateMethods*: Number of private methods in classes calculated by the metrics tool RSM  
*Privateattributes*: Number of private attributes in classes calculated by the metrics tool RSM  
*Physicallinesclass*: Number of physical lines in classes calculated by the metrics tool RSM  
*LOCclass*: Lines of code in classes calculated by the metrics tool RSM  
*eLOCclass*: Effective lines of code in classes calculated by the metrics tool RSM  
*lLOCclass*: Logical lines of code in classes calculated by the metrics tool RSM  
*Cycloclass*: Cyclomatic complexity in classes calculated by the metrics tool RSM  
*InterfaceCompClass*: Interface complexity in classes calculated by the metrics tool RSM  
*TotalParamClass*: Total parameters in classes calculated by the metrics tool RSM  
*TotalReturnClass*: Total returns in classes calculated by the metrics tool RSM  
*CommentsClass*: Comments in classes calculated by the metrics tool RSM  
*BlanksClass*: Number of blank lines in classes calculated by the metrics tool RSM  
*Case*: Number of occurrence of the key word “case” calculated by the metrics tool RSM  
*KWBreak*: Number of occurrence of the key word “break” calculated by the metrics tool RSM  
*LOCh*: Lines of code in .h files calculated by the metrics tool RSM  
*If*: Number of occurrence of the key word “if” calculated by the metrics tool RSM  
*Else*: Number of occurrence of the key word “else” calculated by the metrics tool RSM  
*{}*: Number of lines with key words “{” or “}” calculated by the metrics tool RSM  
*Goto*: Number of occurrence of the key word “goto” calculated by the metrics tool RSM  
*Return*: Number of occurrence of the key word “return” calculated by the metrics tool RSM  
*()*: Number of lines with key words “(” or “)” calculated by the metrics tool RSM  
*Exit*: Number of occurrence of the key word “exit” calculated by the metrics tool RSM  
*\_exit*: Number of occurrence of the key word “\_exit” calculated by the metrics tool RSM  
*Abort*: Number of occurrence of the key word “abort” calculated by the metrics tool RSM  
*eLOCh*: Effective lines of code in .h files calculated by the metrics tool RSM  
*Macro*: Number of occurrence of the key word “macro” calculated by the metrics tool RSM  
*Union*: Number of occurrence of the key word “union” calculated by the metrics tool RSM  
*lLOCh*: Logical lines of code in .h files calculated by the metrics tool RSM  
*Class*: Number of occurrence of the key word “class” calculated by the metrics tool RSM

*Blankh*: Blank lines in .h files calculated by the metrics tool RSM  
*Commenth*: Comments in .h files calculated by the metrics tool RSM  
*Inline*: Inline comments in .h files calculated by the metrics tool RSM  
*TotalLogicalh*: Total logical lines in .h files calculated by the metrics tool RSM  
*Memoryalloc*: Number of occurrence of the key word “alloc” for memory calculated by the metrics tool RSM  
*Memoryfree*: Number of occurrence of the key word “free” for memory calculated by the metrics tool RSM  
*TotalPhysicalh*: Total physical lines in .h files calculated by the metrics tool RSM  
*Memorynew*: Number of occurrence of the key word “new” for memory calculated by the metrics tool RSM  
*Memorydelete*: Number of occurrence of the key word “delete” for memory calculated by the metrics tool RSM  
*Literalstrings*: Number of occurrence of literal strings calculated by the metrics tool RSM  
*Continuation*: Number of occurrence of line continuations calculated by the metrics tool RSM  
*Preprocessor*: Number of preprocessor lines calculated by the metrics tool RSM  
*Include*: Number of “include” calculated by the metrics tool RSM  
*Define*: Number of “defines” calculated by the metrics tool RSM  
*Typedef*: Number of “typedef” calculated by the metrics tool RSM  
*Const*: Number of occurrence of the key word “const” calculated by the metrics tool RSM  
*Enum*: Number of occurrence of the key word “enum” calculated by the metrics tool  
*Do*: Number of occurrence of the key word “do” calculated by the metrics tool RSM  
*While*: Number of occurrence of the key word “while” calculated by the metrics tool RSM  
*Switch*: Number of occurrence of the key word “switch” calculated by the metrics tool RSM  
*Default*: Number of occurrence of the key word “default” calculated by the metrics tool RSM  
*For*: Number of occurrence of the key word “for” calculated by the metrics tool RSM  
*Baseclass*: Number of base classes calculated by the metrics tool RSM  
*Derivedclass*: Number of base classes calculated by the metrics tool RSM  
*Quality*: Number of quality notices by the metrics tool RSM

## **12.2 Development metrics**

*UpdateNotCFiles*: Number of updates to files that are not C source files during development  
*CUpdate*: Number of updates to files that are C source files during development  
*TotalUpdate*: Total number of updates during development  
*NotcAdded*: Number of lines added to files that are not C source files during development  
*CAdded*: Number of lines added to files that are C source files during development  
*Added*: Total number of lines added during development  
*Notcdeleted*: Number of lines deleted from files that are not C source files during development  
*Cdeleted*: Number of lines deleted from files that are C source files during development  
*Deleted*: Total of lines deleted during development



*Modified*: Total number of lines modified during development  
*Difference*: Lines added minus lines deleted during development  
*Notnumauthors*: Number of authors of changes to files that are not C source files during development  
*Cnumauthors*: Number of authors of changes to files that are C source files during development  
*Totalnumauthors*: Total number of authors of changes during development  
*BotHalfnotC*: Number of changes to files that are not C source files by authors in the bottom half in terms of number of changes  
*BotHalfC*: Number of changes to files that are C source files by authors in the bottom half in terms of number of changes  
*BotHalfTotal*: Number of changes to files by authors in the bottom half in terms of number of changes  
*PreBugsAll*: Number of bugs during development for all active releases  
*PreBugsPrev*: Number of bugs during development for the previous release  
*PreBugsCurrent*: Number of bugs during development for the current release  
*PreBugsUnknown*: Number of bugs during development where the release is unspecified  
*Months*: Months between first defect and deployment

### **12.3 Deployment and usage metrics**

*MiscMailings*: Number of messages to the misc mailing list during development.  
*AdvocayMailings*: Number of messages to the advocacy mailing list during development.  
*AnnounceMailings*: Number of messages to the announcement mailing list during development.  
*PortsMailings*: Number of messages to the ports mailing list during development.  
*WWWMailings*: Number of messages to the www mailing list during development.  
*BugsMailings*: Number of messages to the bugs mailing list during development.  
*TechMailings*: Number of messages to the technical mailing list during development.  
*ChangeRequests*: Change requests during development  
*DocBugs*: Document problems during development

### **12.4 Software and hardware configurations metrics**

*AllbugsotherHW*: Number of bugs during development for all active releases against all other HW  
*AllBSDBugsi386HW*: Number of bugs during development for all active releases against x86 HW  
*AllBSDBugssparcHW*: Number of bugs during development for all active releases against sparc HW  
*CurrentBSDBugsotherHW*: Number of bugs during development for the current releases against all other HW  
*CurrentBSDBugsi386HW*: Number of bugs during development for the current releases against x86 HW  
*CurrentBSDBugssparcHW*: Number of bugs during development for the current releases against sparcHW  
*SparcMailings*: Number of messages to the sparc mailing list during development.

## 13 APPENDIX C

In this section, we present the values of the metrics we collected for releases 2.4 – 3.3 (except release 3.2).

<i>Metric</i>	<i>Release 2.4</i>	<i>Release 2.5</i>	<i>Release 2.6</i>	<i>Release 2.7</i>	<i>Release 2.8</i>	<i>Release 2.9</i>	<i>Release 3.0</i>	<i>Release 3.1</i>	<i>Release 3.3</i>
LOC	5170933	5340424	5810112	6228123	6557530	6525985	6700415	6777811	7077564
Statements	2242959	2313823	2534623	2716490	2869484	2844151	2917318	2957934	3072884
Functions	31228	32124	45426	50384	53606	54522	58846	60325	71269
Bandwidth	21810.27	22007.25	24804.13	26279.62	27346.78	26504.45	27611.73	27929.09	28358.65
DeepNest	3791	3845	4214	4415	4783	4759	4799	4805	4921
Cstatements	1957580	2020030	2172704	2329374	2464715	2440893	2503682	2538388	2646759
Lineswithcomments	949837	979421	1030722	1092413	1139499	1116084	1138301	1149453	1172001
InlineComment	127414	130503	135267	143383	150505	144534	143562	144280	141255
Blanklines	502780	518458	553743	588185	627752	619611	630671	639626	672765
Linesforpreprocess	238528	246451	262674	283607	299466	294040	294358	293092	298351
Lineswithcode	2752045	2834971	3035060	3260165	3462274	3434900	3526353	3569082	3735768
Totallines	4317450	4450489	4748625	5082631	5380130	5321745	5447723	5508575	5739232
Statements	1481098	1529889	1636381	1753927	1854569	1833851	1878494	1901778	1976739
Commentchars	26601684	27258556	28695142	30379489	31696338	31096977	31862786	32177394	32912915
Nontext commentchars	3794338	3920893	4114823	4298042	4474054	4394199	4413557	4457190	4581923
Whitespacechars	25331060	26120778	27714517	29514385	31362841	30983857	31627665	31957921	33255514
Preprocessorchars	3972081	4113186	4447194	4788392	5090807	5021928	5052895	5081232	5251247
Statementchars	49662768	51802002	55209112	59422303	63286051	63060836	64991967	65870331	69188053
Totalchars	93774007	97474610	1.04E+08	1.11E+08	1.17E+08	1.15E+08	1.18E+08	1.2E+08	1.23E+08
Tokens	7777859	8052375	8660048	9297996	9899501	9815601	10078167	10184386	10700109
UOpand	1137692	1172554	1260575	1360499	1438961	1423165	1475036	1494784	1552950
UOpator	307249	311451	336342	357822	373384	363261	375449	378823	384400
TPpand	6548261	6793547	7306580	7854043	8379425	8305298	8534896	8630505	9114100
TOpator	8923582	9229145	9902440	10626528	11305316	11219266	11520892	11627458	12208559
PGlength	10431149	10742929	11538010	12464450	13195471	13051408	13510330	13687553	14236826

PGvolume	85299041	87955629	94565522	1.02E+08	1.09E+08	1.08E+08	1.12E+08	1.13E+08	1.18E+08
PGeffort	2.76E+10	2.84E+10	3.04E+10	3.22E+10	3.47E+10	3.45E+10	3.51E+10	3.53E+10	3.75E+10
PGlevel	337.648	337.833	487.197	506.528	516.59	500.347	523.047	530.679	539.454
Files	15201	15382	19602	20573	21222	20603	21355	21534	21952
FunctionsRSM	75357	77110	88297	94591	99114	98755	102610	103962	110971
PhysicalLines	2813041	2905392	3113258	3319232	3508784	3490288	3572826	3618972	3841430
LOC RSM	1980497	2024110	2171088	2327761	2473709	2487433	2549945	2566510	2761209
eLOC	1612641	1649675	1765992	1897487	2017332	2029131	2082612	2095935	2260053
ILOC	1179144	1217191	1302510	1395292	1473806	1465026	1508776	1528757	1626512
Cyclo	588455	602196	649176	692881	731564	731577	750783	759062	808138
InterfaceComp	300550	309034	347075	372984	394057	394352	412458	417986	448741
TotalParams	135851	140274	157792	169912	179862	179652	189346	192032	207299
TotalReturn	164699	168760	189283	203072	214195	214700	223112	225954	241442
CommentsRSM	423701	470961	500102	527040	545822	537990	536058	542920	561798
BlanksRSM	247863	254398	270345	286057	308989	310596	315729	317975	343230
Classes	9395	9589	11102	11916	12637	12866	13315	13438	14010
NestedClasses	658	668	764	813	858	893	914	924	940
TotMeth	10083	10366	12380	14040	14701	14697	15007	15112	14276
PubMeth	9177	9460	11477	12828	13489	13485	13791	13896	13239
Publicattributes	63995	66871	73417	79398	84127	82949	85568	86611	89350
ProtMeth	562	562	592	847	847	847	847	847	665
Protectedattributes	228	228	208	242	242	242	242	242	254
PrivateMethods	358	344	311	365	365	365	369	369	372
Privateattributes	556	527	470	538	540	540	549	549	554
PhysicalLinesclass	213047	220271	243134	261277	285487	279507	288515	294611	305725
LOCclass	199542	206158	229342	246888	267545	261497	270566	276481	286505
eLOCclass	180784	186837	205844	222048	240753	234655	242984	248276	257829
ILOCclass	80839	83567	93342	101159	107192	106166	109008	110060	111065
Cycloclass	3746	3748	4959	4930	4941	4941	4975	4975	5091

InterfaceCompClass	5753	5757	7726	7452	7465	7465	7505	7505	7623
TotalParamClass	2750	2752	3772	3539	3541	3541	3556	3556	3606
TotalReturnClass	3003	3005	3954	3913	3924	3924	3949	3949	4017
CommentsClass	83292	86046	90921	95391	101295	97344	96913	97642	95268
BlanksClass	10576	10950	12003	12646	14092	14101	14651	14822	15220
Case	86495	87874	95354	100532	108255	107551	108750	109352	111174
KWBreak	62977	64334	69255	73905	79511	79162	80596	81209	82944
LOCh	3538393	3655496	3978999	4286572	4519206	4525850	4662087	4708862	4907820
If	307082	314631	337716	361854	381186	380853	391279	395282	413931
Else	68679	70080	75147	80230	84356	84547	85861	86573	90858
{}	458169	469401	513858	546704	578945	577260	589012	594710	619409
Goto	16320	16872	18494	19878	20885	21006	22069	22656	25520
Return	129648	133351	148902	160362	169822	170323	177274	179085	187610
()	637	661	712	744	810	841	930	933	998
Exit	5409	5148	5869	6145	6357	6415	6249	6229	5994
_exit	215	214	229	252	264	277	283	349	384
Abort	2727	2723	4183	4204	4503	4409	4460	4464	4526
eLOCh	1829637	3185434	3464429	3739124	3939451	3947749	4072145	4113219	4287413
Macro	33280	34297	38267	41761	43007	42994	43906	43981	46477
Union	824	842	948	982	1038	1024	1036	1042	1060
ILOCh	1650510	1693127	1836233	1962769	2069547	2053046	2105484	2126094	2205548
Class	9323	9586	10955	11769	12487	12716	13165	13288	13860
Blankh	636499	654447	722201	764330	810248	804424	820743	829481	871179
Commenth	1336561	1376375	1468051	1553551	1614221	1616403	1645532	1659336	1722407
Inline	2616	2606	3253	3319	3402	3358	3199	3168	3899
TotalLogicalh	5511553	5686318	6169251	6604453	6943675	6946677	7128362	7197679	7501406
Memoryalloc	3224	3440	3728	4068	4400	4410	4890	4966	5036
Memoryfree	6877	7221	7447	8062	8762	8950	9897	10031	10359
TotalPhysicalh	5297659	5466867	5936264	6355812	6684439	6665131	6850618	6917857	7217021

Memorynew	24	24	29	40	40	39	42	42	47
Memorydelete	32	33	34	41	41	41	47	47	53
Literalstrings	343973	351656	370490	396418	420529	417220	430159	434449	452821
Continuation	71871	72029	83993	86917	87652	86997	88764	88775	93425
Preprocessor	410007	422562	460623	501378	525893	519503	527585	529786	537162
Include	75799	78453	83084	89320	93797	93087	95817	96946	98763
Define	167446	172507	192976	210315	220493	217765	223503	227903	231861
Typedef	6489	6750	8932	9519	9921	9819	9907	9913	10689
Const	22033	22428	28731	31316	35118	36275	37954	38438	43855
Enum	1356	1371	1733	1806	1894	1909	1960	1967	2063
Do	3857	3901	4353	4536	4684	4560	4759	4828	5213
While	22566	22904	24007	25550	26578	26119	26639	26823	27235
Switch	14086	14419	15594	16716	17620	17540	17967	18144	18318
Default	9334	9512	10366	11083	11803	11788	12069	12201	12485
For	32692	33921	36422	38636	40406	40082	41040	41430	42192
Baseclass	9545	9398	10847	11659	12378	12607	13056	13179	13749
Derivedclass	431	260	403	405	408	408	412	412	416
Quality	5022854	5189601	5537801	6000023	6256205	6283992	6469353	6547285	6850235
UpdateNotCFiles	4301	2367	2395	4785	4441	4661	3451	4786	5413
CUpdate	3847	1797	3592	3872	5311	10327	8699	12549	10972
TotalUpdate	8148	4164	5987	8657	9752	14988	12150	17335	16385
NotcAdded	28801	208302	60902	243847	201325	158178	91162	182179	135368
CAdded	59094	170830	140658	184473	283670	301660	204478	375910	425018
Added	87895	379132	201560	428320	484995	459838	295640	558089	560386
Notcdeleted	14212	91377	44481	149151	80924	86522	57663	170533	73729
Cdeleted	38579	99253	110327	103762	145609	227685	161326	292137	319049
Deleted	52791	190630	154808	252913	226533	314207	218989	462670	392778
Modified	140686	569762	356368	681233	711528	774045	514629	1020759	953164
Difference	35104	188502	46752	175407	258462	145631	76651	95419	167608

Notcnumauthors	30	27	38	44	42	59	53	50	61
Cnumauthors	32	26	36	40	43	53	59	55	65
Totalnumauthors	34	29	42	45	46	61	63	63	66
BotHalfnotC	12	3	11	12	12	23	58	12	28
BotHalfC	22	30	12	12	13	103	24	32	129
BotHalfTotal	34	33	23	24	25	126	82	44	157
PreBugsAll	58	62	71	106	140	123	130	186	145
PreBugsPrev	30	34	36	32	51	57	78	66	64
PreBugsCurrent	34	36	32	51	57	78	66	64	65
PreBugsUnknown	3	2	3	4	6	5	17	63	17
Months	3	2	3	3	4	3	3	3	4
MiscMailings	1437	937	3173	4608	8280	6468	5552	6467	8892
AdvocayMailings	105	24	126	87	65	140	109	82	77
AnnounceMailings	9	5	2	17	9	10	11	11	10
PortsMailings	152	149	664	498	1144	1470	1078	1422	1396
WWWMailings	71	82	178	229	386	380	409	628	398
BugsMailings	180	184	189	307	495	598	557	691	687
SparcMailings	61	46	81	127	120	57	225	348	308
TechMailings	536	332	915	842	1445	1174	944	876	881
AllbugsotherHW	6	6	6	14	3	12	22	29	15
AIIBSD Bugs386HW	48	54	42	74	130	92	93	96	106
AIIBSD BugssparcHW	10	9	8	17	5	2	6	11	14
CurrentBSD BugsotherHW	5	2	2	4	0	5	8	14	6
CurrentBSD Bugs386HW	26	34	25	42	55	72	53	53	49
CurrentBSD BugssparcHW	4	1	5	3	2	2	7	2	11
ChangeRequests	6	8	6	11	18	14	22	25	15

DocBugs	3	12	7	13	21	18	33	45	28
---------	---	----	---	----	----	----	----	----	----