

Seesoft—A Tool For Visualizing Line Oriented Software Statistics

Stephen G. Eick, *Member, IEEE*, Joseph L. Steffen, and Eric E. Sumner, Jr.

Abstract—The Seesoft® software visualization system allows one to analyze up to 50 000 lines of code simultaneously by mapping each line of code into a thin row. The color of each row indicates a statistic of interest, e.g., red rows are those most recently changed, and blue are those least recently changed. Seesoft displays data derived from a variety of sources, such as

- version control systems that track the age, programmer, and purpose of the code, e.g., control ISDN lamps, fix bug in call forwarding;
- static analyses, e.g., locations where functions are called; and
- dynamic analyses, e.g., profiling.

By means of direct manipulation and high interaction graphics, the user can manipulate this reduced representation of the code in order to find interesting patterns. Further insight is obtained by using additional windows to display the actual code. Potential applications for Seesoft include discovery, project management, code tuning, and analysis of development methodologies.

Index Terms—Change management systems, code browsing, interactive graphics, line oriented statistics, scientific visualization.

I. INTRODUCTION

A DIFFICULT problem in software engineering is understanding statistics collected at the source code line level of detail. This class of statistics includes information such as who wrote each line, when it was last changed, whether it fixes a bug or adds new functionality, how it is reached, how often it is executed, and so on. The problem is hard for large systems because of the volume of code. A moderately sized system may have thousands of lines of code and a large system may have millions of lines resulting in a large statistical data set. This paper describes a remarkable visualization technique to analyze line oriented data.

Line level data is available on all large software systems. It comes from version control systems, static analyzers, code profilers, and project management tools. This data, however, is underutilized because it is difficult to analyze. Version control systems such as the Revision Control System (RCS) [1], Source Code Control System (SCCS) [2], Change Management System (CMS) [3], Extended Change Management System (ECMS) [4], or SABLE [5] contain a complete history of the code. For each change to the software they typically capture information such as the affected lines, reason for the change, date, and responsible programmer. Static analyzers such as CIA [6] and cscope [7] capture the definitions of functions, types, macros, external variables, etc., and where

they occur in the code. Profilers such as lcomp [8] perform basic block counting, indicating how often individual lines are executed.

Because of the volume of code it is difficult to gain insight from line oriented statistics or to get a perspective on the whole system. Statistical analysis techniques often involve aggregation. For many purposes, however, there is a need for finer grain detail. In addition, aggregation techniques discard the familiar and rich textual representation of the code. Code browsers, code formatting techniques [9], and version editors [10] are useful, but none of these generalizes to study arbitrary line oriented statistics.

Our approach to studying this class of data is to apply Scientific Visualization techniques [11]. We refer to this as *Software Visualization*. There is a distinguished history of visualization research starting with Tufte's seminal work [12]. Previous visualization work has involved traditional statistical data. Some notable examples include *MACSPIN* [13], scatter plot brushing [14], and dynamic graphical methods for analyzing network traffic [15]. Unfortunately, none of these methods is tailored for studying line oriented software data. We know of no techniques for studying this class of data that takes advantage of the underlying textual representation of software.

This paper describes a new technique for visualization and analysis of source code, and a software tool, Seesoft, embodying the technique. There are four key ideas: reduced representation, coloring by statistic, direct manipulation, and capability to read actual code. The reduced representation is achieved by displaying files as columns and lines of code as thin rows. The color of each row is determined by a statistic associated with the line of code that it represents. In several of our examples the statistic will be the date that the line was created. The visual impression is that of a miniaturized copy of the code with color depicting the age of the code. Then, using direct manipulation and high interaction graphics, a user manipulates the display to find interesting patterns. To display the actual code text the user opens up reading windows and positions virtual magnifying boxes over the reduced representation.

Fig. 1 shows a display of a directory containing 20 source code files containing 9 365 lines of code. The height of each column tells the user how large each file is. Files longer than one column are continued over to the next column. For the display, the line color¹ shows the age of each line using a rainbow color scale with the newest lines in red and the oldest

Manuscript received October 1, 1991; revised August 1, 1992. Recommended by R. Selby and K. Torii.

The authors are with AT&T Bell Laboratories, Naperville, IL 60566. IEEE Log Number 9203763.

¹ In black and white versions of this paper color is to be interpreted as gray level. Red is equivalent to dark grey, green to medium gray, and blue to light gray.

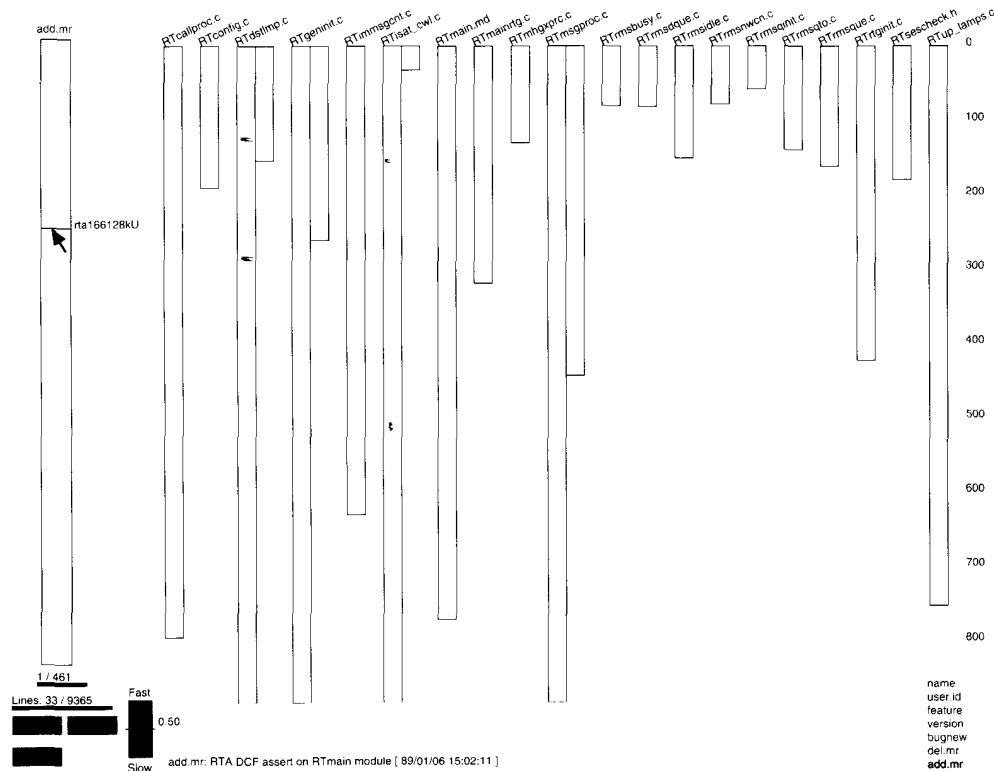


Fig. 2. Example MR activation. As the user positions the mouse over an MR, the lines of code that MR created are activated.

2.1. Screen Description

The Seesoft screen layout consists of a file display, a mouse sensitive color scale, buttons, toggles, and a list of statistic names. Fig. 1 shows that the largest portion of the screen display consists of files shown as columns containing lines of code shown as colored rows. Using a 1280×1024 standard high-resolution monitor we can display about 900 lines of code per column. Files longer than 900 lines wrap and are displayed as multiple columns. For example, file `RTmsgproc.c` in the middle of Fig. 1 has 1 300 lines and is displayed in two columns. The name of each file is printed above it for easy identification. The row representation shows clearly the indentation and length of each line of code. The color of each line is tied to a line oriented statistic. This statistic is highlighted on the list of statistic names in the lower right-hand corner. The rows are just large enough so that block comments, functions, and control structures such as *case* and *if* statements are visible just by their indentation.

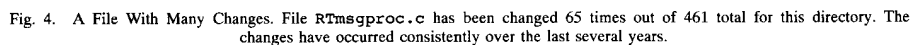
On the left side of the Seesoft display is a mouse sensitive color scale. Each color on the scale represents one value of the statistic associated with each line of code. The statistic might be age, programmer, feature, type of line, number of times the line was executed in a recent test, and so on. We often use the MR (modification request to a version control system) number

for our statistic. The MR number is an interesting statistic because it comes in date order, is the smallest unit of program change, and is associated with programmers, developers, and features. The MR's are displayed sequentially with the newest at the top and oldest at the bottom. Underneath the scale, the number of activated statistic values and the total are shown, 461/461 in Fig. 1, as well as the number of activated code lines and total, 9365/9365. The color scale is a generalization of traditional sliders controlling thresholds. The user may select discontinuous threshold ranges, as well as the traditional continuous ranges that normal sliders may select.

At the bottom of the screen there is space for Seesoft to print the current code line and the statistic value. As the cursor is positioned over any color in the color scale the value of the statistic represented by the color is printed at the bottom of the Seesoft display. In Fig. 1 this is the MR number, abstract (a short description of the purpose of the MR), and date. If the cursor is over a row, Seesoft also prints the line of code associated with that row. We have additional methods of viewing that we describe below.

2.2. Linking Between the Color Bar and Code Lines

Each statistic value is linked to the lines of code having that value through a common color. When the user activates



²A 5ESS expert on this section of the code subsequently told us that these files were added to provide ISDN capability.

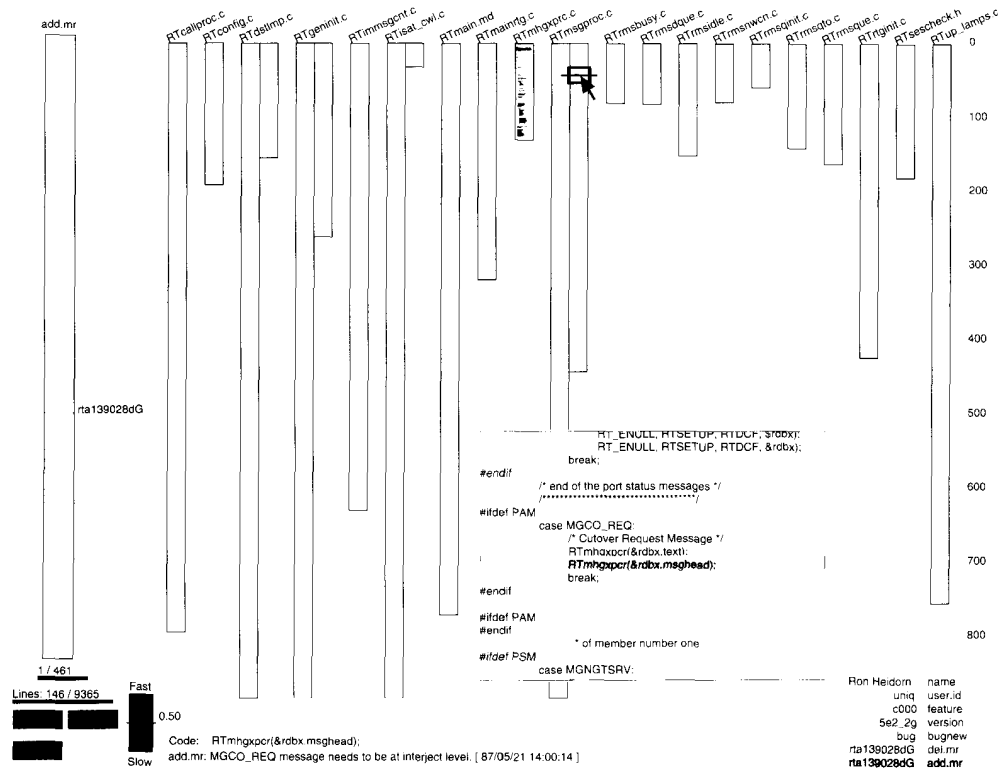


Fig. 5. Modular change. A bug is fixed by inserting a function call in `RTmsgproc.c`, shown in the *Code Reading* window, and putting the code for the function in a new file. The *Code Reading* window may be independently positioned.

been stable since 1985.³ The heavily changed set of files in Fig. 8 are the control flow files in the program. In this directory there is a clear historical work pattern. New functionality is added by creating new functions and then inserting function calls in the main files. Through time there have been several major enhancements that created sets of files.

There are two types of MR's, new feature MR's and bug fixing MR's. Fig. 9 shows the display with the bug MR's in red. Certain files such as the green and light blue files have had few bug fixes. These files were created at one time with a small number of MR's. Other files have multiple bug fixes.

In Fig. 10 the line color is tied to the user id of the programmer writing each line. Files `RTconfig.c`, `RTdstlmp.c`, and `RTup.lamps.c` are written by one individual with a few changes by other people later to fix bugs. Some of the other files with lots of colors have been changed by many different developers. There is a relation between the locations of the bug fixes in Fig. 9 and the number of developers touching the files in Fig. 10. Files touched by many developers have more bug fixes.

³The expert was unaware that these files even existed.

What have we learned in this quick Seesoft session? We know which files are changed most often, the age of the code, and when each file was last changed. We also know that the files in this directory may be clustered into three groups, each group created by a different set of MR's. If it became necessary to divide this directory, files in each cluster could be kept together. We also know where code has been changed recently and that recent MR's have created fewer lines of code than earlier MR's. We also found that certain files have been changed continuously and that the bug fixing MR's are concentrated in these files. These files might be candidates to be restructured or rewritten to reduce maintenance costs.

IV. VISUALIZATION TECHNIQUES

Our approach to visualizing software is to think of source code files and lines as entities in an ordered database. In the preceding examples we display statistics associated with entities obtained from the version control system. For each entity we have a representation, columns, and rows, chosen so that we can view a large volume of data on a single screen. This allows us to gain insight into the overall structure of the database. Database queries are entered and answered visually.

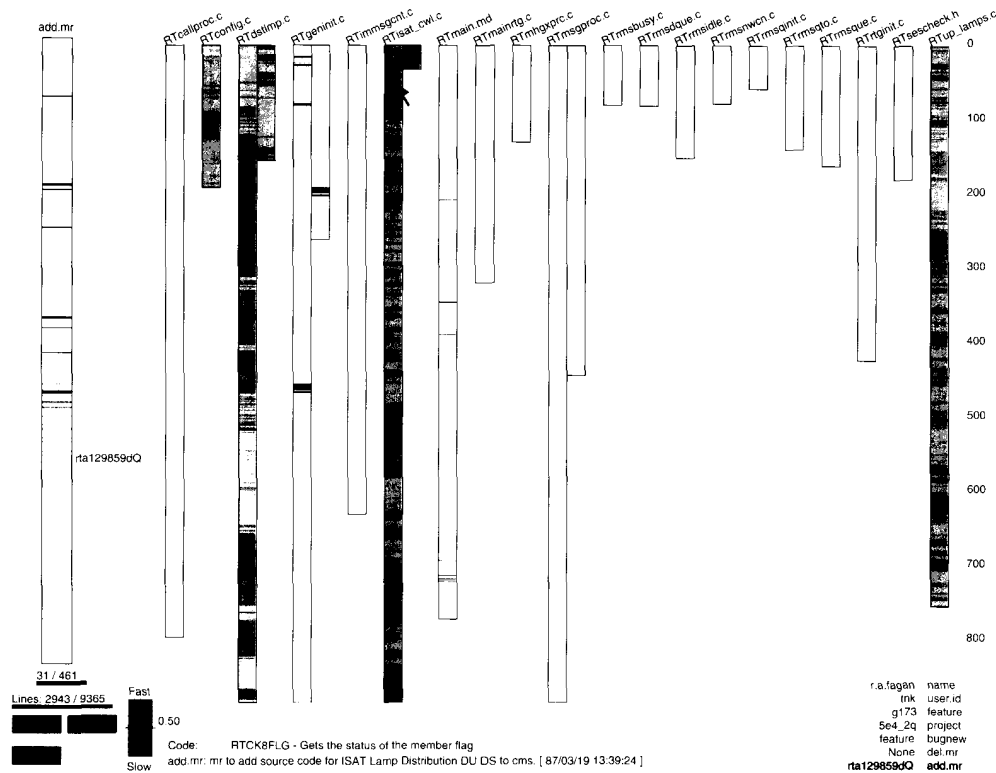


Fig. 6. Groupings of files—a significant new feature. A new feature was put into the code in 1987. This code is primarily in four files and has been stable.

As the user moves the cursor over mouse sensitive portions of the screen he or she is performing a series of database queries. Seesoft then activates the lines of code that resolve the database queries. This approach allows the user to probe the database by moving the mouse around the screen. When he or she discovers an interesting pattern, we provide a mechanism to view the database directly, the reading window in our case, in addition to the reduced representation.

We are currently in the process of obtaining some additional ordered databases. Our approach applies to databases where there is interest in understanding the overall structure and querying the database based on particular attributes. For example, one possible application would be to display a text corpus such as the Bible. Each book could be represented as a column and each verse as a row. A subject index or the age of each verse could be used to color the rows. Another application we are working on is to represent directories as columns and files as rows. This would allow us to visualize even more code on a single display.

The Seesoft user interface employs high interaction graphics and direct manipulation techniques. As the mouse is moved over the display screen entities are automatically activated and deactivated. Since there are no "point and click" delays and no

waiting for screen refreshes, a different style of interaction is possible. Our style of interaction makes it easy for the user to experiment with different activations and to probe the display interactively. For example, with Seesoft it is possible to view each one of several hundred MR's by running the mouse over the color scale. Any unusual MR's will be visually obvious. This would be infeasible if the user were required to "click" on every MR.

V. SEESOFT APPLICATIONS

We envision Seesoft being used in several application areas including

- code discovery,
- new developer training,
- project management,
- quality assurance and system testing,
- software analysis and archeological studies,
- code coverage analysis, and
- code execution optimization.

The code discovery problem is faced by a programmer attempting to change an unfamiliar portion of the source code. Programmers, given requests for additional functionality, must

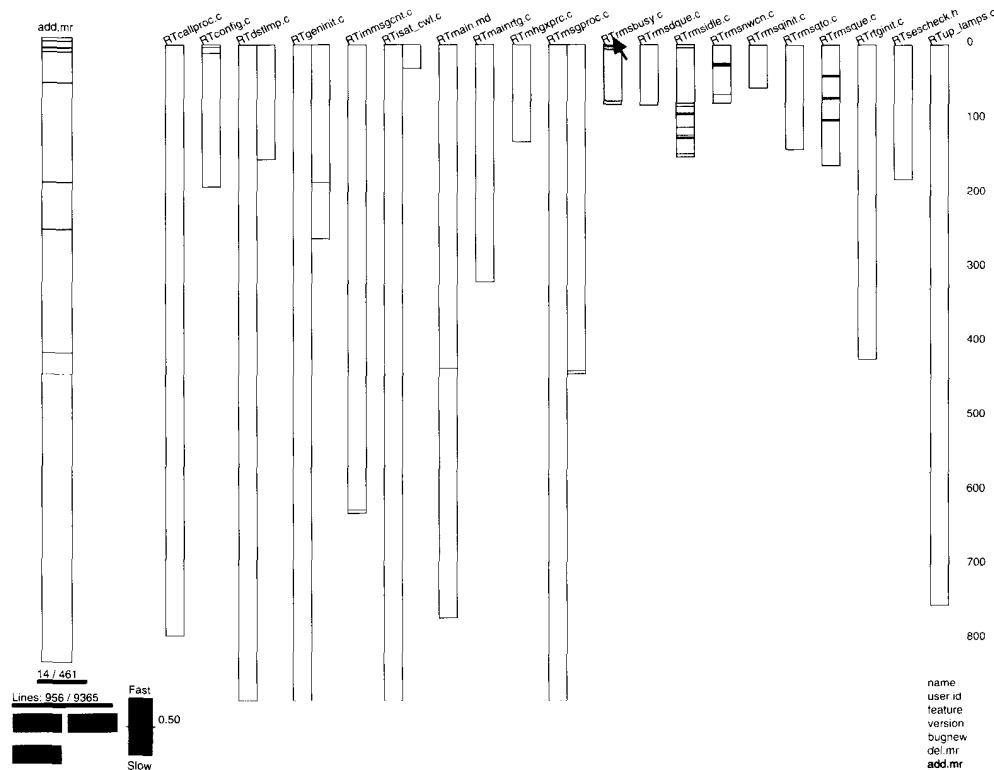


Fig. 7. Groupings of files—stable utility functions. By activating files created by common MR's, we find that there are three different sets of files in this directory. This set is a stable set of utility functions.

study the current code to determine which files contain the existing functionality and which lines to change within these files. This task is often difficult and time consuming. In fact it may take several weeks of detailed study to change a few lines with no unwanted side effects. On large, old projects a significant fraction of a programmer's time is devoted to code discovery. Using Seesoft a programmer can easily determine which lines were created to deliver an existing functionality, the programmers who created those lines, why they were created, and the purpose of nearby lines.

New programmer training is a problem faced by all large software projects. Multiyear projects with large development staffs have considerable staff turnover. Seesoft can ease the programmer training problem by providing new trainees with a global view of the source code. Since Seesoft displays tens of thousands of lines of code simultaneously, new programmers can form a mental picture of the code. Using Seesoft it is easy to answer questions such as:

- How are the files in my program organized?
- When were they last touched?
- Where is the code for this feature?
- What code was written by the person I am replacing?

A class of new programmers might be given Seesoft and access to a code expert. They would use Seesoft to view the code and could immediately ask the expert to explain the interesting things that they discovered.

Project managers monitor a development in order to ensure that the project is on schedule. Using Seesoft a project manager can visually track all source code changes done during the last week or month and can verify that recent changes are consistent with the schedule. In addition, he or she can identify potential trouble spots by the presence of a high level of churn or excessively complex code. A manager can check recent changes in order to trap quick fixes that are likely to cause long term difficulties. He or she can also use Seesoft to identify code that needs to be restructured or rewritten.

Quality assurance inspectors can use Seesoft to determine if new code meets coding specifications and is in the proper files. System testers can determine which regression tests to run by identifying the system functionality embodied in the files that are changed by recent MR's.

Analysts may use Seesoft to understand the effect of various software development environments and processes. For example, an analyst could compare code developed using C

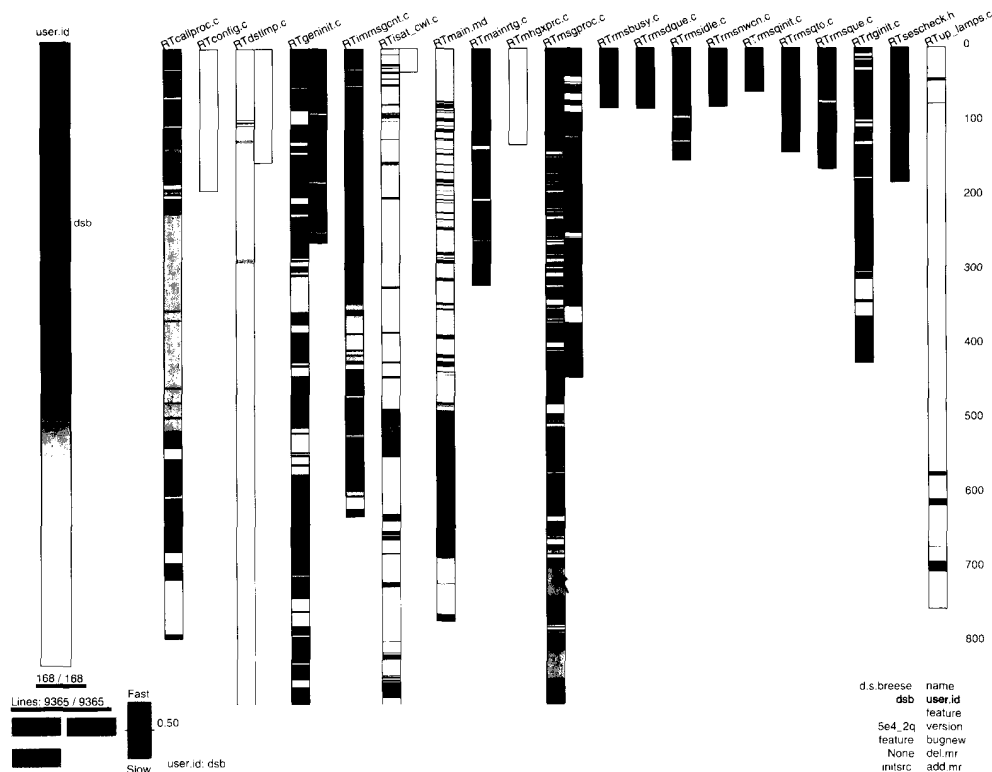


Fig. 9. Locations of bug fixes. MR's for fixing bugs are shown in red. Bug fixes are concentrated in a few of the files.

interested in looking at their own code. One particular expert looked at a recent port and discovered that he had made an error. He had copied a few files from another application and had intended to add a few lines of his own code to each file in order to complete the port. When he colored the copied code in blue, and his own in red, it was immediately obvious that he had failed to add code to one of the files.

VII. IMPLEMENTATION

Seesoft currently runs on Silicon Graphics Iris workstations, although we plan to port it to the X Window System [18]. Seesoft is written in C++ [19] and uses the Silicon Graphics GL graphics library. In total it is about 2 000 lines of C++ code. To deliver real-time user interaction we require the graphics capability to rapidly manipulate the displays, particularly the color map. Seesoft draws each statistic value and associated code lines in its own color. Activating and deactivating is done by manipulating the color map. The colors for the activated statistic values are turned on and for the deactivated values are turned off. Color map manipulation is fast on Iris workstations because it is done in hardware.

Our source code history data came from ECMS. We use the S language for data management and preliminary analysis [20].

To read this data we developed a series of shell scripts to strip unnecessary information and reformat the data for S input. S provides a computational environment, static graphics, and data management that support interactive manipulations. We link Seesoft into the S executive, perform all data manipulation in S, and then launch Seesoft from S.

Silicon Graphics Iris workstations come with 19-in color monitors. Using the column and row representation we find that we can easily understand 20 000 lines of code and can understand 50 000 if we are close to the monitor. In each column we can display about 900 lines and can comfortably fit 25 columns on a single monitor. With more than 50 000 lines displayed the columns become very narrow.

VIII. DISCUSSION AND CONCLUSION

This paper describes a new technique for visualizing line oriented statistics associated with source code and a software tool, Seesoft, embodying the technique. There are four key ideas: reduced representation, coloring by statistic, direct manipulation, and capability to read actual code. The reduced representation is achieved by displaying files as columns and lines of code as thin rows within the columns. The color of each row is determined by a statistic associated

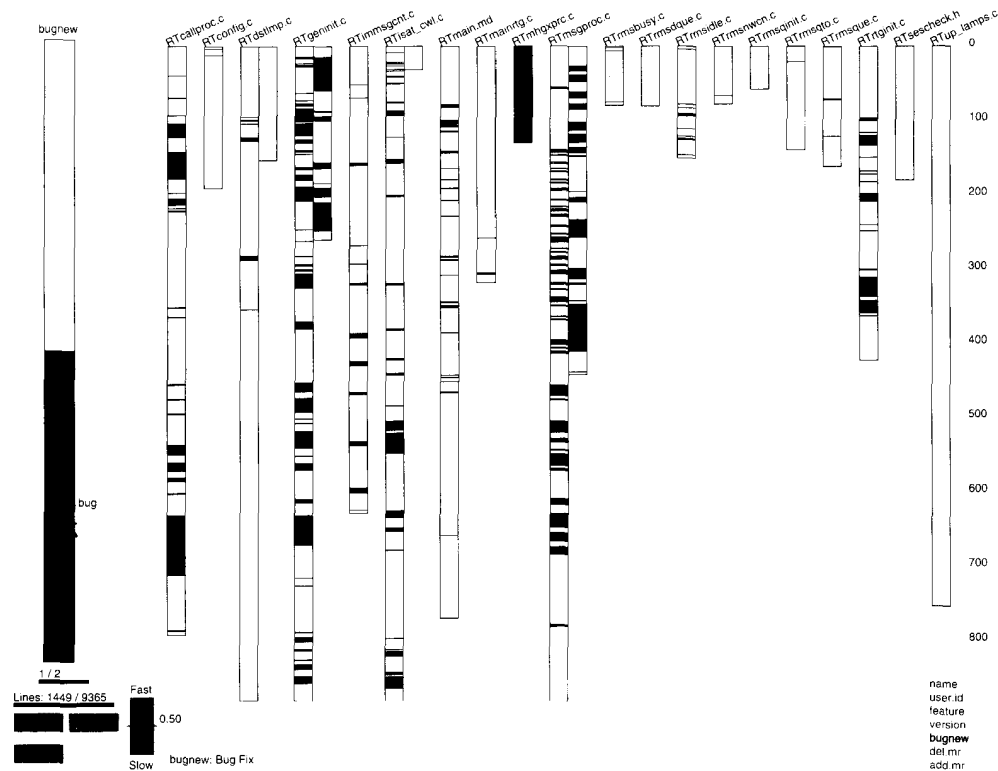


Fig. 10. The color shows the users making changes in this directory. Note that the files with lots of bug fixes are the same files that have been changed by many different developers.

with the line of code that it represents. In our examples we obtained the statistics from a version control system. The visual impression is that of a miniaturized copy of the code with color depicting the spatial distribution of a statistic. Then, using direct manipulation and high interaction graphics, a user manipulates the display to find interesting patterns. To display the actual code text the user opens up reading windows and positions magnifying boxes over the reduced representation.

Besides analyzing source code statistics, our technique has application to any ordered database. Examples include transaction databases, indexed text such as legal writings, a text corpus such as the Bible, and software documentation.

As with any method there are limitations. Our visualization technique provides a qualitative view of the distribution of a statistic in code. As with all graphical methods, the technique is useful for discovering patterns. After the patterns are discovered, hypotheses may be tested by means of standard statistical methods. Currently, Seesoft is unable to display more than 50 000 lines of code simultaneously; however, we are working on other techniques using different abstractions that scale beyond this limit. A key idea in Seesoft is its interactive use of direct manipulation techniques and use of color. It is difficult to describe these in a static monochrome medium such as this

paper.

We developed Seesoft in conjunction with the 5ESS® Telecommunications Switch Project, which includes millions of lines of code developed over 10 years. Because the initial results are so promising we are creating an 85-gigabyte optical disk based archive of the complete code history of 5ESS, including version control data, project management data, and cscope symbol databases derived from monthly snapshots of the code.

ACKNOWLEDGMENT

The authors would like to acknowledge helpful conversations with R. A. Becker, R. Drechsler, G. Nelson, and A. R. Wilks.

REFERENCES

- [1] W. F. Tichy, "RCS—A system for version control," *Software—Practice and Experience*, vol. 15, pp. 637–654, 1985.
- [2] M. J. Rochkind, "The source code control system," *IEEE Trans. Software Engineering*, vol. SE-1, pp. 364–370, 1975.
- [3] B. R. Rowland and R. J. Welsh, "Software development system," *Bell Syst. Tech. J.*, vol. 62, part 2, pp. 275–289, 1983.
- [4] P. A. Tuscany, "Software development environment for large switching projects," in *Proc. Int. Switching Symp.*, pp. 199–214, 1987.

- [5] S. Cichinski and G. S. Fowler, "Product administration through SABLE and NMAKE," *AT&T Tech. J.*, vol. 67, pp. 59-70, 1988.
- [6] Y. F. Chen, "The C program database and its applications," in *Proc. Summer USENIX Conf.*, 1989.
- [7] J. L. Steffen, "Interactive examination of a C program with Cscope," in *USENIX Dallas 1985 Winter Conf. Proc.*, USENIX Association, pp. 170-175, 1985.
- [8] P. J. Weinberger, "Cheap dynamic instruction counting," *AT&T Bell Laboratories Tech. J.*, vol. 63, pp. 1815-26, 1984.
- [9] R. Baecker and A. Marcus, *Human Factors and Typography for More Readable Programs*. Reading, MA: Addison-Wesley, 1990.
- [10] A. A. Pal and M. B. Thompson, "An advanced interface to a switching software version management system," in *Proc. 7th Int. Conf. Software Engineering for Telecommunications Switching Systems*, pp. 110-113, 1989.
- [11] G. M. Nielson, B. Shriver, and L. J. Rosenblum, Eds., *Visualization in Scientific Computing*. Los Alamitos, CA: IEEE Computer Society Press, 1990.
- [12] E. R. Tufte, *The Visual display of Quantitative Information*. Cheshire, CT: Graphics Press, 1983.
- [13] A. W. Donoho, D. L. Donoho, and M. Gasko, *MACSPIN: A Tool for Dynamic Display of Multivariate Data*. Monterey, CA: Wadsworth & Brooks/Cole, 1986.
- [14] R. A. Becker and W. S. Cleveland, "Brushing scatter plots," *Technometrics*, vol. 29, pp. 127-142, 1987.
- [15] R. A. Becker, S. G. Eick, and A. R. Wilks, "Basics of network visualization," *IEEE Computer Graphics and Applications*, vol. 11, pp. 12-14, 1991.
- [16] B. Shneiderman, "Direct manipulation: A step beyond programming languages," *IEEE Computer*, vol. 16, pp. 57-68, 1983.
- [17] R. A. Becker, W. S. Cleveland, and G. Weil, "The use of brushing and rotation for data analysis," pp. 247-275 in *Dynamic Graphics for Statistics*, William S. Cleveland and McGill, Eds. Wadsworth, 1988.
- [18] V. Quercia and T. O'Reilly, "X window system user's guide," O'Reilly & Associates, Inc., Sebastopol, CA, 1988.
- [19] B. Stroustrup, *The C++ Programming Language*. Reading MA: Addison-Wesley, 1987.
- [20] R. A. Becker, J. M. Chambers, and A. R. Wilks, *The New S Language*. Pacific Grove, CA: Wadsworth & Brooks/Cole, 1988.



Stephen G. Eick (M'87) received the B.A. degree from Kalamazoo College in 1980, the M.A. degree in mathematics from the University of Wisconsin, Madison, in 1981, and the Ph.D. degree in statistics from the University of Minnesota in 1985.

He joined AT&T Bell Laboratories in 1985, where he has been involved in statistical graphics and network research. He has developed techniques to visualize networks, network performance models, and network simulations. As a member of the Software Production Research Department, he is applying visualization techniques to understand software.

Dr. Eick is a member of the ACM.



Joe Steffen received the B.S. degree in electrical engineering from Purdue University and the M.S. degree in computer science from the Illinois Institute of Technology.

He has done considerable work on software tools, including writing cscope and ctrace, which are part of UNIX System V, and adding run-time checking of array subscripts and pointer bounds to the portable C compiler. His current research interests are configuration management for large software systems and software tools.

Mr. Steffen is a member of the ACM.



Eric E. Sumner, Jr. received the A.B. and Ph.D. degrees in engineering science from Harvard University in 1980 and 1984, respectively.

He joined AT&T Bell Laboratories in 1984 and is currently Head of the Software Production Research Department, which was formed in 1990 at Indian Hill, the Illinois complex that houses many large software developments, including the SESS® switch. Prior to assuming his current position, he played backgammon professionally, developed models of composite materials, built choice

models for AT&T network equipment, and led the development of a tool for computer-aided engineering of underwater surveillance systems.