

Error Localization During Software Maintenance: Generating Hierarchical System Descriptions from the Source Code Alone

Richard W. Selby† and Victor R. Basili‡

†Department of Information and Computer Science, University of California, Irvine, California 92717¹

‡Department of Computer Science, University of Maryland, College Park, Maryland 20742²

Abstract

One central feature of the structure of a software system is the coupling among its components (e.g., subsystems, modules) and the cohesion within them. The purpose of this study is to quantify ratios of coupling and cohesion and use them in the generation of hierarchical system descriptions. The ability of the hierarchical descriptions to localize errors by identifying error-prone system structure is evaluated using actual error data. Measures of data interaction, called data bindings, are used as the basis for calculating software coupling and cohesion. A 135,000 source line system from a production environment has been selected for empirical analysis. Software error data was collected from high-level system design through system test and from some field operation of the system. A set of five tools is applied to calculate the data bindings automatically, and cluster analysis is used to determine a hierarchical description of each of the system's 77 subsystems. An analysis of variance model is used to characterize subsystems and individual routines that had either many/few errors or high/low error correction effort. The empirical results support the effectiveness of the approach for localizing errors. The approach is especially useful during software maintenance since the tools require only the source code for automatic generation of a hierarchical system description.

1 Introduction

Several researchers have proposed methods for relating the structure of a software system to its quality (e.g., [BE82] [HK81] [Eme84]). One pivotal step in assessing the structure of a software system is characterizing its coupling and cohesion. Intuitively, the *cohesion* in a software system is the amount of interaction *within* pieces (e.g., subsystems, modules) of a system. Correspondingly, *coupling* in a software system is the amount of interaction *across* pieces of a system. Cohesion may sometimes be referred to as "strength."

¹This work was supported in part by the National Science Foundation under grant CCR-8704311 with cooperation from the Defense Advanced Research Projects Agency (ARPA Order 6108, Program Code 7T10), by the National Science Foundation under grant DCR-8521398, and by IBM under the Shared University Research (S.U.R.) Program.

²This work was supported in part by IBM under the Shared University Research (S.U.R.) Program.

Various interpretations for coupling and cohesion have been proposed [SMC74]. In this paper, we present an empirical study that investigates hierarchical software system descriptions that are based on measures of cohesion and coupling. The study evaluates the effectiveness of the hierarchical descriptions in identifying error-prone system structure. Our measurement of cohesion and coupling is based on intra-system interaction in terms of *software data bindings* [BT75] [HB85]. Our measurement of error-proneness is based on software error data collected from high-level system design through system test; some error data from system operation are also included.

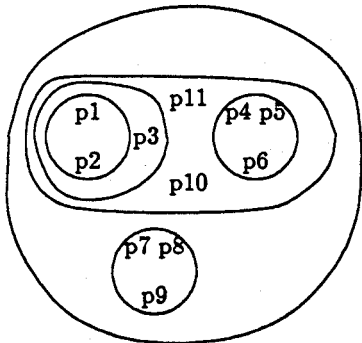
Section 2 discusses the software project selected. The data bindings software analysis and supporting tools are described in Section 3. The data analysis appears in Section 4. Section 5 presents the interpretations and conclusions.

2 Selected Software Project

The software project selected for study is the next release of an internal software library tool. The previous system release contains approximately 100,000 source lines. The production of the next release requires the development or modification of approximately 40,000 source lines. Hence, the total size of the next system release is approximately 135,000 source lines. The system is written in four languages: a high-level programming language similar to PL/I, a language for operating system executives, a user-interface specification language, and an assembly language. The static source code metrics discussed later, including the data bindings analysis, pertain to only the system portion written in the high-level source language. This portion constitutes approximately 70% of the lines in the system and the vast majority of the system logic and intra-system interactions. Project duration, including system and field test, spanned approximately 16 months and maximum staffing included 23 persons. The error data were collected and analyzed at the same time the project took place. An important goal was to minimize the impact of the data collection process on the developers. See [SB88] for a description of the data collection process.

There are 163 source code files in the system containing a total of 451 source code *routines*. A routine is a main program, procedure, or function. The number of routines

Figure 1: Example hierarchical cluster based on software data bindings. Procedures and functions are denoted by p_i , and clusters are denoted by circles. The smaller clusters are relatively tighter (and form earlier), while the larger clusters are relatively looser (and form later). The clusters define a system hierarchy in the form of a tree: the smaller clusters at the leaf nodes and the largest cluster at the root node.



per source code file varies from 1 to 21. On the average, there are 2.8 routines per source code file. There are 77 executable features in the system, referred to as *subsystems* in the paper. These subsystems can be thought of as groups of routines collected together to form functional features of the overall system. The number of source files linked together to form a subsystem varies from 1 to 82. On the average, 26.3 source files are linked together into a subsystem. The same source file is bound into 12.4 different subsystems on the average. Subsystems averaged 19,000 source lines in size, including comments.

3 Cluster Analysis Using Data Bindings

One primary goal for this study was to investigate the relationship of "software data bindings" to software errors [HB85]. "Data bindings" are measures that capture the data interaction across portions of a software system. The theoretical background for the measures is described in [HB85]. Earlier studies have revealed insights about the usefulness of data bindings in the characterization of software systems and their errors [BT75] [HB85].

A *data binding* is defined as an ordered triple (p,x,q) , where p and q are procedures and x is a variable within the static scope of both p and q , and p assigns a value to x and q references x . Data bindings reflect the possibility of a data interaction between two components, based upon p , q , and x . Data bindings count occurrences where there may be a flow of information from p to q via the variable x . The possible orders of execution for p and q are not considered. That is, there may be other factors (e.g., control flow conditions) which would prevent such communication.

In [HB85], this type of data bindings is called "actual data bindings." There are types of data bindings that measure

stronger levels of interaction. Actual data bindings were chosen for this study since they seem to offer an adequate measure of similarity while not requiring complex data flow analysis that stronger levels need. Essentially, we are erring in the direction of safety (as done, for example, by code optimizers) by assuming that procedures may influence one another unless we can show otherwise.

First, we calculated the data bindings in the system. Then, we applied the statistical technique of clustering [Eve80] to the data bindings information to produce a hierarchical description for the software system (see Figure 1). The clustering takes place in a bottom-up manner. The process iteratively creates larger and larger clusters, until all the elements have collapsed into a single cluster. The elements in the clusters are the procedures and functions in the system. The elements with the greatest interaction, in terms of data bindings, cluster together. For more details about the calculation of data bindings and hierarchical clusters in this study, see [SB88]. The technique of clustering has been applied previously to partition a large system into subsystems in [BE82]. Hierarchical clusters have been formally defined in [JS71].

Applicability During Software Maintenance

A set of five software tools was developed to calculate these hierarchical, data bindings clusters and applied to the 77 subsystems in the selected project. For a description of the tools, see [SB88]. The trees provide a form of system documentation — they give a hierarchical view of the subsystems with respect to data usage (see Figure 1). The tools require only the source code for automatic generation of a hierarchical system description. Therefore, the approach is especially useful during software maintenance since — all too often — all that remains of a system is the source code itself.

4 Data Analysis

4.1 Terminology

Throughout the analysis and interpretation, we use the terms *subsystems* and *routines* as follows:

- Routine — A routine is a main program, procedure, or function. There are a total of 451 source code routines in the system.
- Subsystem — A subsystem is a large set of routines that are linked together to form an executable system feature. There are 77 executable features in the system. They average 19,000 source lines in size.

A routine is linked into 12.4 subsystems on the average. Therefore, the total size of the whole system is not $77 \times 19,000 = 1,463,000$ source lines; the total size is approximately 135,000 source lines. See Section 2 for further de-

Figure 2: Distribution of errors and error correction effort by subsystem coupling/strength ratios.

Subsystem coupling/strength	Errors				Error correction hours			
	per KLOC		Total		per KLOC		Total	
	Mean	Std	Mean	Std	Mean	Std	Mean	Std
High	1.54	3.95	0.44	0.99	2.80	7.53	0.88	2.69
Low	0.31	1.16	0.15	0.52	0.91	4.51	0.42	2.39
Overall	1.28	3.58	0.38	0.92	2.39	7.03	0.78	2.63

Figure 3: Distribution of errors and error correction effort by subsystem size.

Subsystem size	Errors				Error correction hours			
	per KLOC		Total		per KLOC		Total	
	Mean	Std	Mean	Std	Mean	Std	Mean	Std
Large	1.52	3.94	0.43	0.98	2.77	7.44	0.86	2.61
Small	0.35	1.22	0.17	0.58	0.98	4.96	0.49	2.71
Overall	1.28	3.58	0.38	0.92	2.39	7.03	0.78	2.63

scription of the subsystems and routines in the software system.

We used the analysis tools mentioned in Section 3 to produce hierarchical descriptions for each of the 77 subsystems (see Figure 1). The hierarchical descriptions are rooted, connected trees that indicate the internal subsystem structure. Each routine in a subsystem occurs as a leaf node in the tree exactly once. Subtrees indicate groupings of routines that form natural *clusters* based on the data bindings criteria. There is a one-to-one correspondence between subtrees and clusters. A cluster can contain either routines or other clusters. In other words, the root node of a subtree can have as its children either leaf nodes (i.e., routines) or the root node of another subtree (i.e., a subset of its own routines that form a smaller cluster).

In the software system being analyzed, a routine may be linked into more than one subsystem. Each of the 77 subsystems has a separate hierarchical description. Therefore, a routine appears in the hierarchical description of each subsystem into which it is linked. A routine may cluster with different sets of routines in different subsystems.

Associated with each cluster in a subsystem is a number that reflects the nature of the binding of the routines in the cluster. This number is interpreted as the following ratio:

$$\frac{\text{the coupling of the cluster with other clusters in the subsystem}}{\text{the internal strength of the cluster}}$$

That is, the number captures the coupling/strength ratio for a cluster of routines within a subsystem. Software engineering principles generally suggest that it is desirable to have low coupling and high strength, which in this context means a low coupling/strength ratio [SMC74].

The data bindings analysis produced 77 trees corresponding to the subsystems. We calculated three different mea-

sures based on the clusters resulting from the data bindings analysis. For each routine occurrence, we calculated:

- Routine coupling/strength ratio — The coupling/strength ratio of the first cluster to form that included the routine as a member. This metric is intended to capture the relationship of a routine to other routines in a subsystem in terms of coupling and strength.
- Routine location in subsystem's data binding tree — The depth in the tree of the first subtree (i.e., cluster) to form that included the routine as a member. More precisely, it is the depth in the tree of the root of that subtree. This metric is intended to characterize the location of a routine in a data binding tree. This location information is useful to know when data binding trees are used as an alternate form of system documentation.

For each subsystem, we calculated:

- Subsystem coupling/strength ratio — The median of the coupling/strength ratios for the clusters within the subsystem. We use a non-parametric statistic here, i.e., a median, because the coupling/strength ratios are relative measures. This metric is intended to characterize the overall coupling and strength within a subsystem.

4.2 Data Analysis Method

An analysis of variance model was used to characterize subsystems and routines that had either many/few errors or high/low development effort spent in error correction.

4.2.1 Independent Variables

The analysis of variance model [Sch59] considered numerous factors simultaneously: subsystem size (above/below median); subsystem coupling/strength ratio (above/below

Figure 4: Distribution of errors and error correction effort across subsystem coupling/strength ratios and subsystem size.

Subsystem coupling/strength	Subsystem size	Errors				Error correction hours			
		per KLOC		Total		per KLOC		Total	
		Mean	Std	Mean	Std	Mean	Std	Mean	Std
High	Large	1.66	4.12	0.46	1.02	2.99	7.71	0.92	2.66
	Small	0.45	1.41	0.21	0.66	1.11	5.31	0.56	2.93
Low	Large	0.36	1.27	0.15	0.52	0.91	4.11	0.39	2.07
	Small	0.28	1.09	0.15	0.52	0.90	4.75	0.44	2.57
Overall		1.28	3.58	0.38	0.92	2.39	7.03	0.78	2.63

Figure 5: Distribution of errors and error correction effort by routine coupling/strength ratios.

Routine coupling/strength	Errors				Error correction hours			
	per KLOC		Total		per KLOC		Total	
	Mean	Std	Mean	Std	Mean	Std	Mean	Std
4_Highest	2.27	4.58	0.59	1.04	5.86	10.98	1.94	4.20
3_Higher	1.15	3.13	0.34	0.74	2.19	6.84	0.72	2.54
2_Lower	1.45	4.19	0.44	1.18	1.57	4.27	0.49	1.61
1_Lowest	0.28	1.11	0.15	0.49	0.21	1.09	0.06	0.29
Overall	1.28	3.58	0.38	0.92	2.39	7.03	0.78	2.63

median); individual subsystem's attributes (77 levels); routine size (above/below median); routine coupling/strength ratio (split into four quartiles); routine location in subsystem's data binding tree (split into four quartiles); and two-way interactions. When defining the levels for some of the factors, non-parametric statistics (e.g., medians, quartiles) were used since the coupling/strength ratios are relative measures and the data bindings trees have different overall depths. Subsystem size and routine size are included as factors in the analysis because earlier analyses have indicated a relationship between size and software effort and error data (e.g., [Boe81] [BSP83]). For a more complete description of the factors and their levels, see [SB88].

4.2.2 Dependent Variables

There were four dependent variables examined with the analysis of variance model.

1. Total errors — The total number of inspection, Trouble Report (TR), System Trouble Report (STR), and Error Summary Worksheet (ESW) errors in a routine³
2. Total errors per KLOC — The total number of inspection, TR, STR, and ESW errors in a routine per 1000 lines of source code
3. Error correction effort — The total amount of effort (in hours) spent correcting TR and ESW errors in a

³Inspections were held during the high-level and low-level design phases and after the completion of unit testing. Error Summary Worksheet (ESW) errors were recorded during the coding, unit testing, and integration testing phases. System Trouble Report (STR) errors were recorded during system testing. Trouble Report (TR) errors were reported against working, released code during and after field testing.

routine

4. Error correction effort per KLOC — The total amount of effort (in hours) spent correcting TR and ESW errors in a routine per 1000 lines of source code

In general, the discussion will focus on the errors per KLOC and the error correction effort per KLOC measures of the routines as opposed to the absolute numbers. This factors out possible underlying correlations between source lines and number of errors or amount of error correction effort. The statistics for all four measures are reported, however. The discussion will tend to highlight results that demonstrated a statistically significant difference, as opposed to those where there was no statistical difference. All results discussed are statistically significant at least at the $\alpha < .05$ level.

4.3 Characterization of High-Error and Low-Error Subsystems

In the source code portions of the system (see Section 2), there was a total of 299 distinct errors recorded from inspections, error summary worksheets (ESW's), system trouble reports (STR's), and trouble reports (TR's). Data on the effort required for error correction were available for 204 distinct errors recorded on ESW's and TR's. In the subsequent figures, all inspection, ESW, STR, and TR errors are counted equally.

Figures 2, 3, and 4 present the errors and error correction effort in the routines in subsystems with: different coupling/strength ratios, different sizes, and different combinations of coupling/strength ratio and size, respectively.

Figure 6: Distribution of errors and error correction effort by routine size.

Routine size	Errors				Error correction hours			
	per KLOC		Total		per KLOC		Total	
	Mean	Std	Mean	Std	Mean	Std	Mean	Std
Large	1.19	2.54	0.47	0.99	3.22	8.72	1.20	3.42
Small	1.39	4.55	0.26	0.80	1.36	3.81	0.26	0.71
Overall	1.28	3.58	0.38	0.92	2.39	7.03	0.78	2.63

Figure 7: Distribution of errors and error correction effort by routine location in data binding tree.

Routine tree location	Errors				Error correction hours			
	per KLOC		Total		per KLOC		Total	
	Mean	Std	Mean	Std	Mean	Std	Mean	Std
4_Root	0.88	2.82	0.30	0.77	1.30	4.82	0.37	1.59
3_Shallow	1.78	4.44	0.51	1.12	3.55	8.88	1.19	3.36
2_Deep	0.96	2.48	0.27	0.63	2.51	7.39	0.83	2.82
1_Deep	1.28	3.73	0.38	0.96	1.76	5.08	0.57	1.95
Overall	1.28	3.58	0.38	0.92	2.39	7.03	0.78	2.63

Graphical plots of the data are presented in [SB88].

Summary of Results

1. Large subsystems with high coupling/strength ratios had routines with the most errors per KLOC.
2. Large subsystems with high coupling/strength ratios had routines with six times as many errors per KLOC than did small subsystems with low coupling/strength ratios.
3. Large subsystems with high coupling/strength ratios had routines with ten times as many unit and integration test (ESW⁴) errors per KLOC than did small subsystems with low coupling/strength ratios.
4. Large subsystems with high coupling/strength ratios had routines with eight times as much error correction effort per KLOC from unit and integration test (ESW) errors than did small subsystems with low coupling/strength ratios.
5. Large subsystems had routines with more errors per KLOC than did small subsystems.

4.4 Characterization of High-Error and Low-Error Routines

As mentioned in Section 4.3 there were 299 distinct errors, counting all inspection, ESW, STR, and TR errors equally; 204 of them had data on error correction effort. Figures 5, 6, and 7 present the errors and error correction effort in the routines with: different coupling/strength ratios, different sizes, and different data binding tree locations, respectively. Various graphical plots of the data are presented in [SB88].

⁴Errors during the coding and unit and integration testing phases were reported on error summary worksheets (ESW's).

Summary of Results

1. The routines with the highest coupling/strength ratios (4_HIGHEST) had the most errors per KLOC and the most error correction effort per KLOC.
2. The routines with the lowest coupling/strength ratios (1_LOWEST) had the fewest errors per KLOC and the least error correction effort per KLOC.
3. The routines with the highest coupling/strength ratios had over eight times as many errors per KLOC than did routines with the lowest coupling/strength ratios.
4. The routines with the highest coupling/strength ratios had over 27 times as much error correction effort per KLOC than did routines with the lowest coupling/strength ratios.
5. Routines in data binding tree location region 3_SHALLOWER had more errors per KLOC and more error correction effort per KLOC than did routines in the other tree regions.
6. Routines that had both the highest coupling/strength ratios (4_HIGHEST) and a location in the "central portion" of the data binding tree (3_SHALLOWER or 2_DEEPER) had the most error correction effort per KLOC.
7. Small routines had more unit and integration test (ESW) errors per KLOC than did large routines.
8. Large routines had more error correction effort per KLOC than did small routines when either all errors or just unit and integration test (ESW) errors were considered.
9. Large routines tended to have a higher average amount of correction effort per error for unit and integration test (ESW) errors than did small routines.

5 Interpretations and Conclusions

In this study, we have merged two goals:

- To collect and analyze data from an ongoing software project without negatively impacting the software developers; and
- To investigate hierarchical system descriptions based on the software engineering principles of coupling and strength (or cohesion) and their relationship to software errors and error correction effort.

This study highlights and empirically supports several software engineering principles. The interpretations span several areas: coupling/strength, system structure, and size.

Coupling/Strength

Low coupling/strength ratios are desirable (i.e., these results empirically support the software engineering principle of desiring low coupling and high strength).

- Routines with the lowest coupling/strength ratios had 8.1 times fewer errors per KLOC than routines with the highest coupling/strength ratios and errors were 27.9 times less costly to fix.
- Large subsystems with high coupling/strength ratios had routines with 4.6 times more errors per KLOC than did the other categories of subsystems.

System Structure Hierarchy: Data Bindings View

The structure of the system at the highest level, i.e., initial stages of problem decomposition, and lowest level, e.g., formulation of abstract data types, appear to be better understood than the intermediate levels of abstraction and specification.

- The errors were 50% less costly to fix in routines at the shallowest (4_ROOT) and deepest (1_DEEPEST) levels of the data bindings view of the system structure hierarchy than at the middle levels, and there were 21% fewer errors per KLOC.

Size

Subsystem size seems to be at least as important, if not more important, than routine size. Hence, maybe the software community has been worrying about the wrong issue.

- Smaller subsystems had routines with 4.3 times fewer errors per KLOC than did larger subsystems.
- Smaller routines had a slightly higher average of errors per KLOC than did larger routines, although the difference was not statistically significant. When just unit and integration test errors are considered, however, smaller routines had significantly more errors per KLOC than did larger routines. Overall, errors in smaller routines were 2.4 times less expensive to fix.

6 Acknowledgement

The authors are very grateful to several persons on the selected software project for their assistance and support in this research. Their names cannot be mentioned because of a non-disclosure agreement. The authors appreciate the assistance of D. Hutchens in developing the data bindings analysis tools and S. Wilkin in collecting the data.

References

- [BE82] L.A. Belady and C.J. Evangelisti. System partitioning and its measure. *Journal of Systems and Software*, 2(1):23-29, February 1982.
- [Boe81] B. W. Boehm. *Software Engineering Economics*. Prentice-Hall, Englewood Cliffs, NJ, 1981.
- [BSP83] V. R. Basili, R. W. Selby, and T. Y. Phillips. Metric analysis and data validation across fortran projects. *IEEE Transactions on Software Engineering*, SE-9(6):652-663, Nov. 1983.
- [BT75] V. R. Basili and A. J. Turner. Iterative enhancement: a practical technique for software development. *IEEE Transactions on Software Engineering*, SE-1(4), Dec. 1975.
- [Eme84] T. Emerson. A discriminant metric for module cohesion. In *Proc. Seventh Intl. Conf. Software Engr.*, pages 294-303, Orlando, FL, 1984.
- [Eve80] B. S. Everitt. *Cluster Analysis, 2nd ed.* Heineman Educational Books Ltd., London, 1980.
- [HB85] D. H. Hutchens and V. R. Basili. System structure analysis: clustering with data bindings. *IEEE Trans. Soft. Engr.*, SE-11(8), Aug. 1985.
- [HK81] S. Henry and D. Kafura. Software quality metrics based on interconnectivity. *Journal of Systems and Software*, 2(2):121-131, 1981.
- [JS71] N. Jardine and R. Sibson. *Mathematical Taxonomy*. John Wiley and Sons, New York, 1971.
- [SB88] Richard W. Selby and Victor R. Basili. *Analyzing Error-Prone System Coupling and Cohesion*. Technical Report, Dept. of Computer Science, University of California, Irvine, 1988.
- [Sch59] H. Scheffe. *The Analysis of Variance*. John Wiley & Sons, New York, 1959.
- [SMC74] W. P. Stevens, G. J. Myers, and L. L. Constantine. Structural design. *IBM Systems Journal*, 13(2):115-139, 1974.